# A neural basis of probabilistic computation in visual cortex

Edgar Y. Walker[1,2†], R. James Cotton[1,2,3†], Wei Ji Ma[4‡], Andreas S. Tolias[1,2,5*‡]

[1]*Center for Neuroscience and Artificial Intelligence, Baylor College of Medicine, TX, USA*

[2]*Department of Neuroscience, Baylor College of Medicine, TX, USA*

[3]*Now at: Shirley Ryan AbilityLab, IL, USA*

[4]*Center for Neural Science and Department of Psychology, New York University, NY, USA.*

[5]*Department of Electrical and Computer Engineering, Rice University, TX, USA.*

†‡ These authors contributed equally to this work

∗ Corresponding author

**For more than a century, Bayesian-inspired models have been used to explain human and animal behavior, suggesting that organisms represent the uncertainty associated with sensory variables. Nevertheless, the neural code of uncertainty remains elusive. A central hypothesis is that uncertainty is encoded in the population activity of cortical neurons in the form of likelihood functions. Here, we studied the neural code of uncertainty by simultaneously recording population activity in the visual cortex in primates during a visual categorization task for which trial-to-trial uncertainty about stimulus orientation was relevant for the animal's decision. We decoded the likelihood function from the trial-to-trial population activity and found that it predicted the monkey's decisions better than using only a decoded point**

**estimate of the orientation. Critically, this remained true even when we conditioned on the stimulus including its contrast, suggesting that random fluctuations in neural activity firing rates drive behaviorally meaningful variations in the likelihood function. Our results establish the role of population-encoded likelihood functions in mediating behavior and offer potential neural underpinnings for Bayesian models of perception.**

When making perceptual decisions, organisms often benefit from representing uncertainty about sensory variables. More specifically, the theory that the brain performs Bayesian inference—which has roots in the work of Laplace[1] and von Helmholtz[2]—has been widely used to explain human and animal perception[3–6]. At its center lies the assumption that the brain maintains a statistical model of the world and when confronted with incomplete and imperfect information, makes inferences by computing probability distributions over task-relevant variables. In spite of the prevalence of Bayesian theories in neuroscience, evidence to support them stems primarily from behavioral studies (e.g.[7,8]), and how probability distributions or uncertainty are encoded in the brain remains unclear.

According to the probabilistic population coding (PPC) hypothesis[9,10], inference in the brain is performed by inverting a generative model of neural population activity. Specifically, according to PPC, a neural population encodes sensory uncertainty in the form of the sensory likelihood function—the probability of observing a given pattern of neural activity across hypothesized stimulus values[9,11,12]. The form of the likelihood function is inherited from the probability distribution describing neural variability ("noise") for a given stimulus. A sensory likelihood function is typi-

41 cally unimodal[13,14], and its width could in principle serve as a measure of the sensory uncertainty

42 about the stimulus. Whether the brain uses this particular uncertainty quantity in its decisions is

43 unknown.

44 Testing this hypothesis thoroughly is not straightforward. Experiments to test for any population-

45 level neural code for trial-by-trial uncertainty must: (1) Use a behavioral task in which uncertainty

46 information, and not just a point estimate, is relevant for the perceptual decision, (2) Record simul-

47 taneously from a population of neurons, and (3) Show that the stimulus-conditioned fluctuations in

48 uncertainty decoded from population responses—for us, the sensory likelihood function—mediate

49 the behavioral outcomes (perceptual decisions). Despite previous efforts[10,15,16], these three criteria

50 have not yet been met simultaneously, and therefore, the population neural code of uncertainty still

51 remains unknown.

52 To illustrate the importance of the third criterion, consider a typical perceptual decision-making

53 task (Fig. 1a) where the subject views a stimulus $s$, which elicits a cortical population response $\mathbf{r}$,

54 for example in V1. Here by a stimulus, we refer collectively to all aspects of a visual stimulus such

55 as its contrast and orientation. Stimulus information is eventually relayed to decision-making areas

56 (e.g. prefrontal cortex), leading the animal to make a classification decision $\hat{C}$. The likelihood

57 function $\mathcal{L}$ is decoded from the recorded population activities $\mathbf{r}$. Because variation in the stimulus

58 (e.g. orientation or contrast) across trials can drive variation both in the decoded likelihood function

59 and in the animal's decision, one may find a significant correlation between $\mathcal{L}$ and $\hat{C}$, even if the

60 likelihood estimated from the recorded population $\mathbf{r}$ does not mediate the decision (Fig. 1c). When

3

61    the stimulus is fixed, random fluctuations in the population response $\mathbf{r}$ can still result in variations

62    in $\mathcal{L}$. If the likelihood truly mediates the decision, we expect that such variation in $\mathcal{L}$ would account

63    for variation in $\hat{C}$. Therefore, to demonstrate that the likelihood $\mathcal{L}$ mediates decision $\hat{C}$ (Fig. 1d),

64    it is imperative to show a correlation between $\mathcal{L}$ and $\hat{C}$ conditioned on the stimulus $s$ (Fig. 1d,e).

65    In this work, we meet these requirements and provide the first evidence that in perceptual decision-

66    making, cortical populations encode and utilize trial-by-trial sensory uncertainty information in

67    the form of likelihood functions. We performed simultaneous V1 cortical population recordings as

68    monkeys performed a visual classification task in which the trial-by-trial uncertainty information

69    is beneficial to the animal[17]. To decode the trial-by-trial likelihood functions from the V1 popula-

70    tion responses, we developed a novel technique based on deep learning[18,19]. This method extends

71    beyond previous approaches that used strong parametric assumptions about the probability dis-

72    tribution of population activity[10,14,16]. These assumptions were theoretically convenient[9,20–23] but

73    limited the generality of those approaches in decoding the likelihood functions. Finally, we per-

74    formed all analyses conditioned on the contrast—an overt driver of uncertainty—and performed

75    further orientation-conditioned analyses to isolate the effect of random fluctuations in the decoded

76    likelihood function on behavior.

## Results

78    **Behavioral task**   Two Rhesus macaques (*Macacca mulatta*) were trained on an orientation classifi-

79    cation task designed such that the optimal performance required the use of trial-by-trial uncertainty.
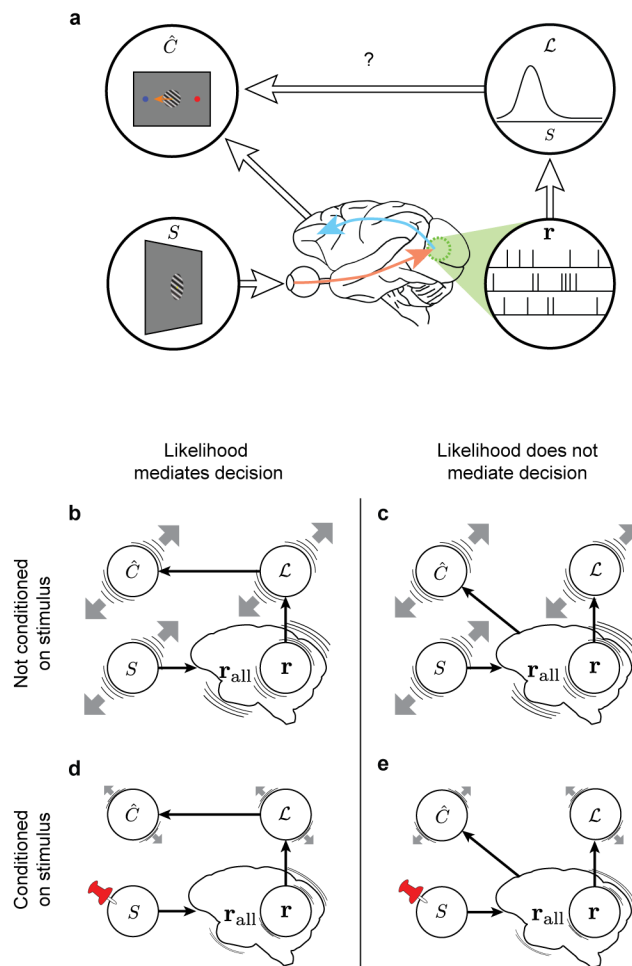
**Figure 1:** Conceptual overview of the decoding model. **a**, Information flow between stimulus $s$, recorded cortical population response $\mathbf{r}$, responses of all recorded and unrecorded neurons $\mathbf{r}_{\text{all}}$, decoded likelihood function $\mathcal{L}$, and subject's decision $\hat{C}$, as the subject performs a visual classification task. **b-e**, Possible relationship between variables in the model indicated by black arrows. We consider two scenarios: **b, d** the likelihood function mediates the decision, **c, e** the likelihood function does not mediate the decision. The gray arrow represents the trial-by-trial fluctuations in the subject's decisions $\hat{C}$ as predicted by the variable. **b, c**, When not conditioning on the stimulus, the stimulus can drive correlation among all variables, making it difficult to distinguish the two scenarios. **d, e**, When conditioning on the stimulus, we expect correlation between $\hat{C}$ and $\mathcal{L}$ only when $\mathcal{L}$ mediates the decision, allowing us to distinguish the two scenarios.

80  On each trial, one of two stimulus classes ($C = 1$ or $C = 2$) was chosen at random with equal

81  probability. Each class was defined by a Gaussian probability distribution over the orientation.

82  The two distributions shared the same mean but had different standard deviations (Fig. 2a). An

83  orientation was drawn from the distribution belonging to the selected class, and a drifting grating

84  stimulus with that orientation was then presented to the animal (Fig. 2b). In a given recording

85  session, at least three distinct contrasts were selected at the beginning of the session, and on each

86  trial, one of these values was randomly selected.

87  In our previous study[17], we designed this task so that an optimal Bayesian observer would in-

88  corporate the trial-by-trial sensory uncertainty about stimulus orientation in making classification

89  decisions. Indeed, both humans and monkeys decisions seemed to utilize trial-by-trial uncertainty

90  about the stimulus orientation. In the current study, one of the two monkeys (Monkey L) was the

91  same monkey that participated in the previous study and thus has been shown to have learned the

92  task well. The second monkey (Monkey T) was naïve to the task, but learned to perform equally

93  well, closely matching the performance of Monkey L (Fig. 2c), with psychometric curves display-

94  ing a strong dependence on both contrast and orientation (Fig. 2d,e).

95  In our analyses, we grouped the trials with the same contrast within the same session and refer to

96  such a group as a "contrast-session".

97  **Decoding a cortical population representation of uncertainty** Each monkey was implanted

98  with a chronic multi-electrode (Utah) array in the parafoveal primary visual cortex (V1) to record

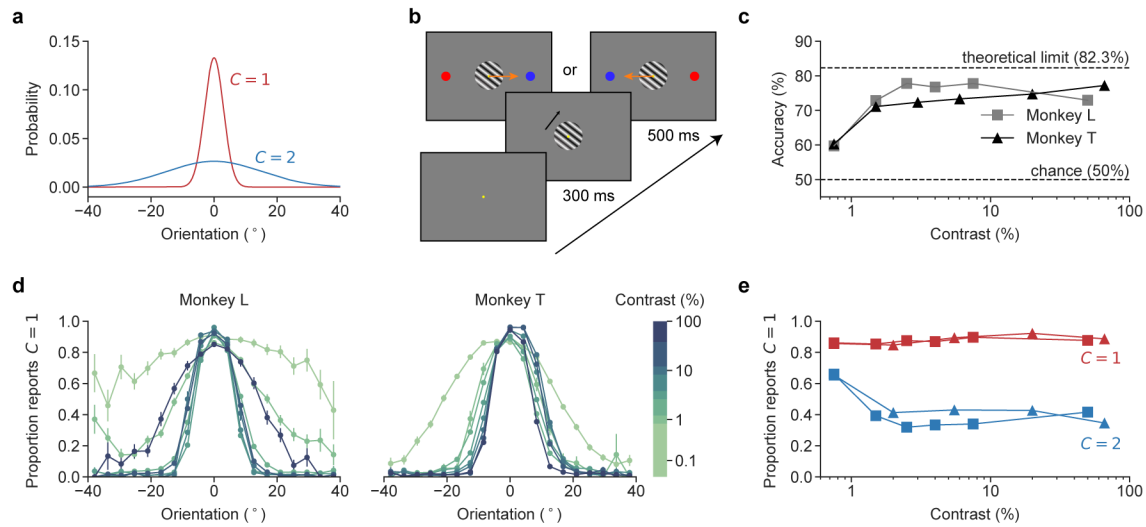99  the simultaneous cortical population activity as they performed the orientation classification task

**Figure 2:** Behavioral task. **a,** The stimulus orientation distributions for the two classes. The two distributions shared the same mean ($\mu = 0°$) but differed in their standard deviations ($\sigma_1 = 3°$ and $\sigma_2 = 15°$). **b,** Time course of a single trial. The subject fixated onto the fixation target for 300 ms before a drifting grating stimulus was shown. After 500 ms of stimulus presentation, the subject broke fixation and saccaded to one of the two colored targets to indicate their class decision (color matches class color in **a**). The left-right configuration of the colored targets were chosen at random for each trial. **c,** Performance of the two monkeys on the task across stimulus contrast. "Theoretical limit" corresponds to the performance of an ideal observer with no observation noise. **d,** Psychometric curves. Each curve shows the proportion of trials on which the monkey reported $C = 1$ as a function of stimulus orientation, computed from all trials within a single contrast bin. All data points are means and error bars indicate standard error of the means. **e,** Class-conditioned responses. For each subject, the proportions of $C = 1$ reports is shown across contrasts, conditioned on the ground-truth class: $C = 1$ (red) and $C = 2$ (blue). The symbols have the same meaning as in **c**.

100  (Fig. 3a). A total of 61 and 71 sessions were analyzed from Monkey L and Monkey T for a to-

101  tal of 110,695 and 192,631 trials, respectively. In each recording session, up to 96 channels were

102  recorded. On each trial and for each channel, we computed the total number of spikes that occurred

103  during the 500 ms of stimulus presentation preceding the decision-making cue (Fig. 3a), yielding

104  a vector of population responses $\mathbf{r}$ used in the subsequent analyses (Fig. 3b).

105  Existing computational methods for decoding the trial-by-trial likelihood function from the cor-

106  tical population activities are typically based on making strong parametric assumptions about the

107  stimulus conditioned distribution of the population response (i.e. generative model of the pop-

108  ulation response). For example, population responses to a stimulus can be modeled as an inde-

109  pendent Poisson distribution, allowing each recorded unit to be characterized by a simple tuning

110  curve (which may be further parameterized)[14, 16, 20–23]. While this simplifying assumption allows

111  the trial-by-trial likelihood function to be computed straightforwardly, it disregards any potential

112  correlations among the units in population responses (i.e. noise correlations and internal brain state

113  fluctuations[24–29]), and can lead to biased estimates of the likelihood encoded by the cortical popu-

114  lation. While more generic parametric models—such as Poisson-like distributions—of population

115  distribution have been proposed[9, 10, 17, 30, 31], they still impose restrictive assumptions.

116  We devised a technique based on deep learning to decode the trial-by-trial likelihood function from

117  the V1 population response. We trained a fully connected deep neural network (DNN)[19] to predict

118  the per-trial likelihood function $\mathcal{L}(\theta) \equiv p(\mathbf{r}|\theta)$ over stimulus orientation $\theta$ from the vectorized

119  population response $\mathbf{r}$ (Fig. 3c; for details on the network architecture and training objective, refer
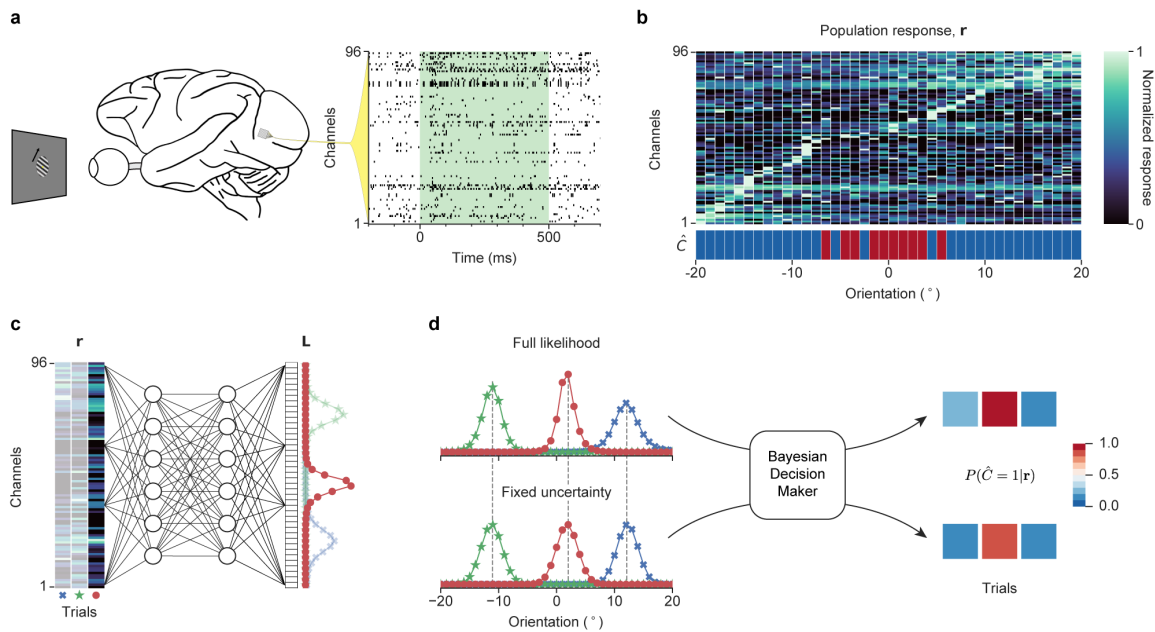
8

**Figure 3:** Encoding and decoding of the stimulus orientation. **a**, An example of 96 channels spike traces from a single trial from Monkey T. The vector of spike counts, **r**, was accumulated over the pre-saccade stimulus presentation period (time 0-500 ms, green shade). **b**, The population response for the selected trials from a single contrast-session (Monkey T, 64% contrast). Each column is a population response **r** on a trial randomly drawn from the trials falling into a specific orientation bin. Each row is a response from a single channel. For visibility, the channel's responses are normalized to the maximum response across all trials. The channels were sorted by the preferred orientation of the channel. The subject's class decision $\hat{C}$ for each trial is depicted by the colored patch: red and blue for $\hat{C} = 1$ and $\hat{C} = 2$, respectively. **c**, A schematic of a deep neural network trained to decode the vectorized likelihood function **L** from **r**. All likelihood functions are area-normalized. **d**, Two decision models $M$. In the *Full-Likelihood Model*, the decoded likelihood function **L** was used without modification to predict the probability that the monkey will report $\hat{C} = 1$. In the *Fixed-Uncertainty Model*, the likelihood function was approximated by a fixed width Gaussian centered at the mean of the normalized likelihood (dashed line), where the fixed width was fitted separately for each contrast-session. For both models, their likelihood functions were fed into a parameterized Bayesian decision maker to yield the decision prediction $p(\hat{C} = 1|\mathbf{r})$, where parameters of the decision maker were fitted to the subject's classification decisions, separately for each contrast-session.

9

120 to Methods). A separate network was trained for each contrast-session and no behavioral data was

121 utilized in training the DNN. By using a DNN to directly learn the likelihood decoder, we avoid

122 making any parametric assumption on the population response distribution. DNNs can also learn to

123 decode likelihood functions for neural populations with known parametric generative models such

124 as Poisson-like distributions, and therefore constitute a strictly more flexible likelihood decoding

125 method. On simulated population responses, trained DNNs could well recover the ground-truth

126 likelihood functions (Extended Data Fig. 1; refer to Methods for the simulation details).

127 The likelihood functions decoded by the DNNs exhibit the expected dependencies on the overt

128 drivers of uncertainty, contrast (Fig. 4a-c): the width of the likelihood function is higher at lower

129 contrast (Fig. 4d).

130 **Decision-making model** As we varied the stimulus contrast from trial to trial, the expected un-

131 certainty about the stimulus orientation varied, and one would expect the monkeys to represent

132 and make use of their trial-by-trial sensory uncertainty in making decisions. However, we make a

133 much stronger claim here: even at a fixed contrast, because of random fluctuations in the popula-

134 tion response[32,33], we predict (1) the uncertainty encoded in the population, that is, the likelihood

135 function, to still fluctuate from trial to trial, and (2) the effect of such fluctuations to manifest in the

136 monkey's decisions (Fig. 1d) on a trial-by-trial basis. We tested this prediction by fitting, separately

137 for each contrast-session, the following two models and comparing their performance in predicting

138 the monkey's decision (Fig. 3d): (1) a Full-Likelihood Model, which utilizes the trial-by-trial un-

139 certainty information decoded from the population response in the form of the likelihood function,
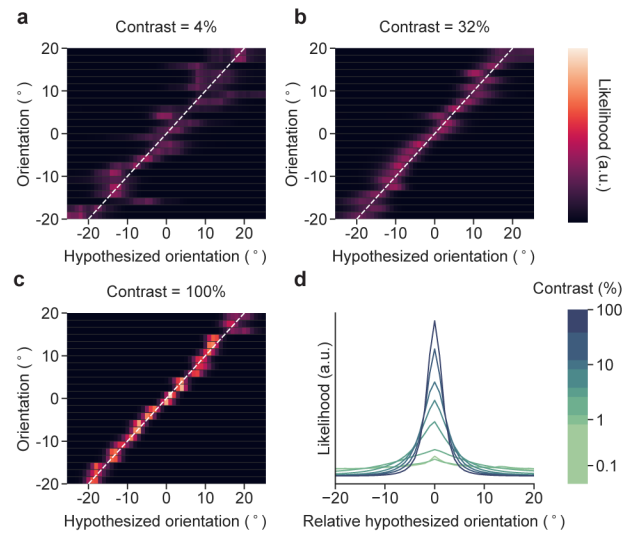
10

**Figure 4:** Likelihood functions decoded by the trained neural networks. **a-c,** Example decoded likelihood functions from three contrast-sessions from Monkey T. Each row represents the decoded likelihood function over the hypothesized orientation for a randomly selected trial within the specific orientation bin. All likelihood functions are area-normalized. Brighter colors correspond to higher values of the likelihood function. **d,** Average likelihood function by contrast. On each trial, the likelihood function was shifted such that the mean orientation of the normalized likelihood function occured at $0°$. The centered likelihood functions were then averaged across all trials within the same contrast bin.

140 and (2) a Fixed-Uncertainty Model, which approximates the trial-by-trial likelihood function by a

141 Gaussian function centered at the mean of the normalized likelihood function, with a fixed width

142 that is fitted separately for each contrast-session.

143 In both models, the (original or Gaussian approximated) likelihood functions were fed into the

144 Bayesian decision maker to yield trial-by-trial prediction of the subject's decision in the form of

145 $p(\hat{C}|\mathbf{r}, M)$, or the likelihood of subject's decisions $\hat{C}$ conditioned on the population response $\mathbf{r}$

146 and the model $M$. The Bayesian decision maker computed the posterior probability of each class

147 and used these to produce a stochastic decision. The means of the class distributions assumed by

148 the observer, the class priors, the lapse rate, and a parameter to adjust the exact decision-making

149 strategy were used as free parameters (refer to Methods for details). The model parameters were

150 fitted by maximizing the total log likelihood over all trials in a contrast-session $\sum_i \log p(\hat{C}_i|\mathbf{r}_i, M)$.

151 The fitness of the models was assessed through cross-validation, and we reported mean and total

152 log likelihood of the models across all trials in the test set.

153 Both models incorporated a trial-by-trial point estimate of the stimulus orientation (the mean of the

154 normalized likelihood function) and only differed in whether they contain additional uncertainty

155 information about the stimulus orientation carried by the trial-by-trial fluctuations in the shape of

156 the likelihood function. We use the term "shape" to refer to all aspects of the likelihood function

157 besides its mean, including its width. If the fluctuations in the shape of the likelihood function truly

158 captured the fluctuations in the sensory uncertainty as represented and utilized by the animal, one

159 would expect the Full-Likelihood Model to yield better trial-by-trial predictions of the monkey's

160 decisions than the Fixed-Uncertainty Model.

161 We observed that the Full-Likelihood Model performed well above chance in predicting the mon-

162 key's behavior across all contrasts (Extended Data Fig. 2), reaching up to 90% accuracy. In mean

163 log likelihood, the Full-Likelihood Model consistently outperformed the Fixed-Uncertainty Model

164 across contrasts and for both monkeys (Fig. 5a,b; trial log likelihood differences between the Full-

165 Likelihood and Fixed-Uncertainty Model: Monkey L: paired t-test, $t(110694) = 25.32$, $p < 0.001$,

166 $\delta_{\text{total}} = 1975.90$ and Monkey T: $t(192610) = 19.23$, $p < 0.001$, $\delta_{\text{total}} = 1611.50$; $\delta_{\text{total}}$ is the total

167 log likelihood difference across all trials). This shows that the trial-by-trial fluctuations in the shape

168 of the likelihood function are informative about the monkey's trial-by-trial decisions, demonstrat-

169 ing that decision-relevant sensory uncertainty information is contained in population responses

170 that can be captured in the shape of the likelihood function. This in turn strongly supports the

171 hypothesis that visual cortex encodes stimulus uncertainty in the form of a likelihood function on

172 the trial-by-trial basis.

173 We next asked how meaningful our effect sizes (model performance differences) are. To answer

174 this question, we simulated the monkey's responses across all trials and contrast-sessions taking

175 the trained Full-Likelihood Model to be the ground truth, and then retrained the Full-Likelihood

176 Model and the Fixed Uncertainty Model from scratch on the simulated data. This approach yields a

177 theoretical upper bound on the observable difference between the two models if the Full-Likelihood

178 Model was the true model of the monkeys' decision-making process. The model performance

179 differences between the Full-Likelihood Model and the Fixed-Uncertainty Model (1975.90 and

13

180 1611.50 total log likelihood differences across all trials for Monkey L and T, respectively) were

181 very close in magnitude to the expected total upper bound log likelihood differences of $2368.25 \pm$

182 $131.45$ and $2355.22 \pm 144.48$ based on the simulations (representing mean $\pm$ standard deviation

183 across 5 repetitions of simulation for Monkey L and T, respectively) (Extended Data Fig. 3). This

184 suggests that our effect sizes are meaningful and that the Full-Likelihood Model is a reasonable

185 approximate description of the monkey's true decision-making process.

186 We next isolated the contribution to the perceptual decision of trial-by-trial fluctuations in the

187 shape of the likelihood function from the contribution of the stimulus orientation, by condition-

188 ing the analysis on the stimulus orientation (Fig. 1d). Specifically, we shuffled the shapes of the

189 decoded likelihood functions across trials within the same orientation bin, separately for each

190 contrast-session. This shuffling preserved the expected uncertainty at each stimulus orientation—

191 the dependence of the likelihood shape on the stimulus orientation—and the trial-by-trial correla-

192 tion between the mean of the likelihood function and the subject's perceptual decision (Fig. 5c),

193 while removing the trial-by-trial correlation between the shape of the likelihood function and the

194 behavioral decision conditioned on the stimulus orientation.

195 By design, the Fixed-Uncertainty Model should perform identically on the original and the shuffled

196 data. If the Full-Likelihood Model outperformed the Fixed-Uncertainty Model simply because it

197 captured spurious correlations between the stimulus orientation and the shape (including the width)

198 of the likelihood function, the performance difference between the two models should remain the

199 same when trained and tested on the shuffled data. However, if the difference in performance
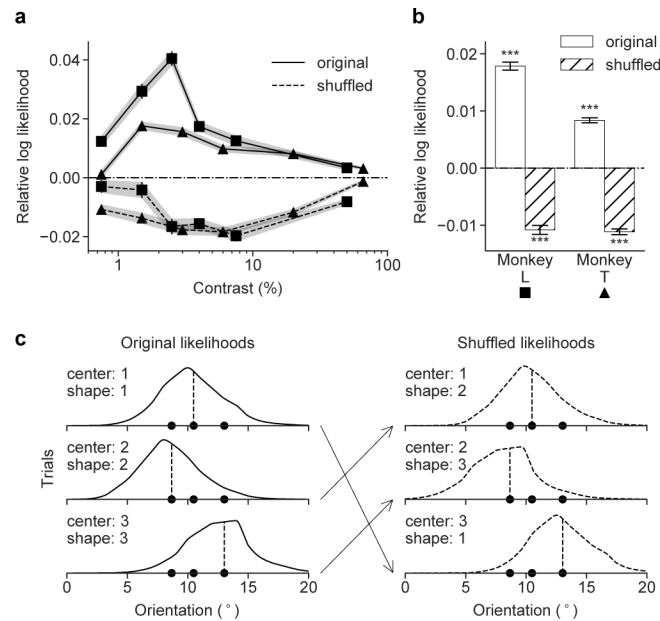
14

**Figure 5:** Model performance. **a**, Average trial-by-trial performance of the Full-Likelihood Model relative to the Fixed-Uncertainty Model across contrasts, measured as the average trial difference in the log likelihood. The results for the original (unshuffled) and the shuffled data are shown in solid and dashed lines, respectively. The squares and triangles mark Monkey L and T, respectively. **b**, Relative model performance summarized across all contrasts. Results are shown for each monkey and for original (solid) and shuffled (hatched) data. The difference between the Full-Likelihood and Fixed-Uncertainty Models was significant with $p < 0.005$ (***) for both monkeys, and on both the original and the shuffled data. Furthermore, the difference between the Full-Likelihood Model on the original and the shuffled data was significant ($p < 0.005$ for both monkeys). For **a** and **b**, all data points are means, and error bar/shaded area indicate standard error of the means. **c**, Shuffling scheme for three example trials drawn from the same stimulus orientation bin. Shuffling maintains the means but swaps the shapes of the likelihood functions.

200 existed because the trial-by-trial fluctuations in the likelihood shape captured the fluctuations in

201 the decision-relevant sensory uncertainty as we hypothesized, one would expect this difference to

202 disappear on the shuffled data.

203 Indeed, the shuffling of the likelihood function shapes abolished the improvement in prediction

204 performance that the Full-Likelihood Model had over the Fixed-Uncertainty Model. In fact, the

205 Full-Likelihood Model consistently underperformed the Fixed-Uncertainty Model on the shuffled

206 data (Fig. 5a,b; trial log likelihood difference between the Full-Likelihood Model and the Fixed-

207 Uncertainty Model on the shuffled data: Monkey L: paired t-test $t(110694) = -13.33$, $p < 0.001$,

208 $\delta_{\text{total}} = -11667.74$ and Monkey T: $t(192610) = -21.19$, $p < 0.001$, $\delta_{\text{total}} = -2097.06$; $\delta_{\text{total}}$ is

209 the total log likelihood difference across all trials). Therefore, there were significant performance

210 differences in Full-Likelihood Model between the unshuffled and shuffled data (trial log likelihood

211 difference: Monkey L: paired t-test $t(110694) = 32.68$, $p < 0.001$, $\delta_{\text{total}} = 3142.64$ and Monkey

212 T: $t(192610) = 34.81$, $p < 0.001$, $\delta_{\text{total}} = 3708.56$).

213 To confirm our effect sizes were appropriate, we again compared these values to those obtained

214 from simulations in which we took the Full-Likelihood Model to be the ground truth (Extended

215 Data Fig. 3). The simulations yielded total log likelihood differences of the Full-Likelihood Model

216 between the unshuffled and shuffled data of $3675.30 \pm 248.42$ (Monkey L) and $4068.98 \pm 173.71$

217 (Monkey T) (mean $\pm$ standard deviation across 5 repetitions), similar in magnitude to the observed

218 values.

219 Taken together, the shuffling analyses show that the better prediction performance of the Full-

16

220 Likelihood Model is not due to the confound between the stimulus and the likelihood shape. We

221 conclude that the trial-by-trial likelihood function decoded from the population represents behav-

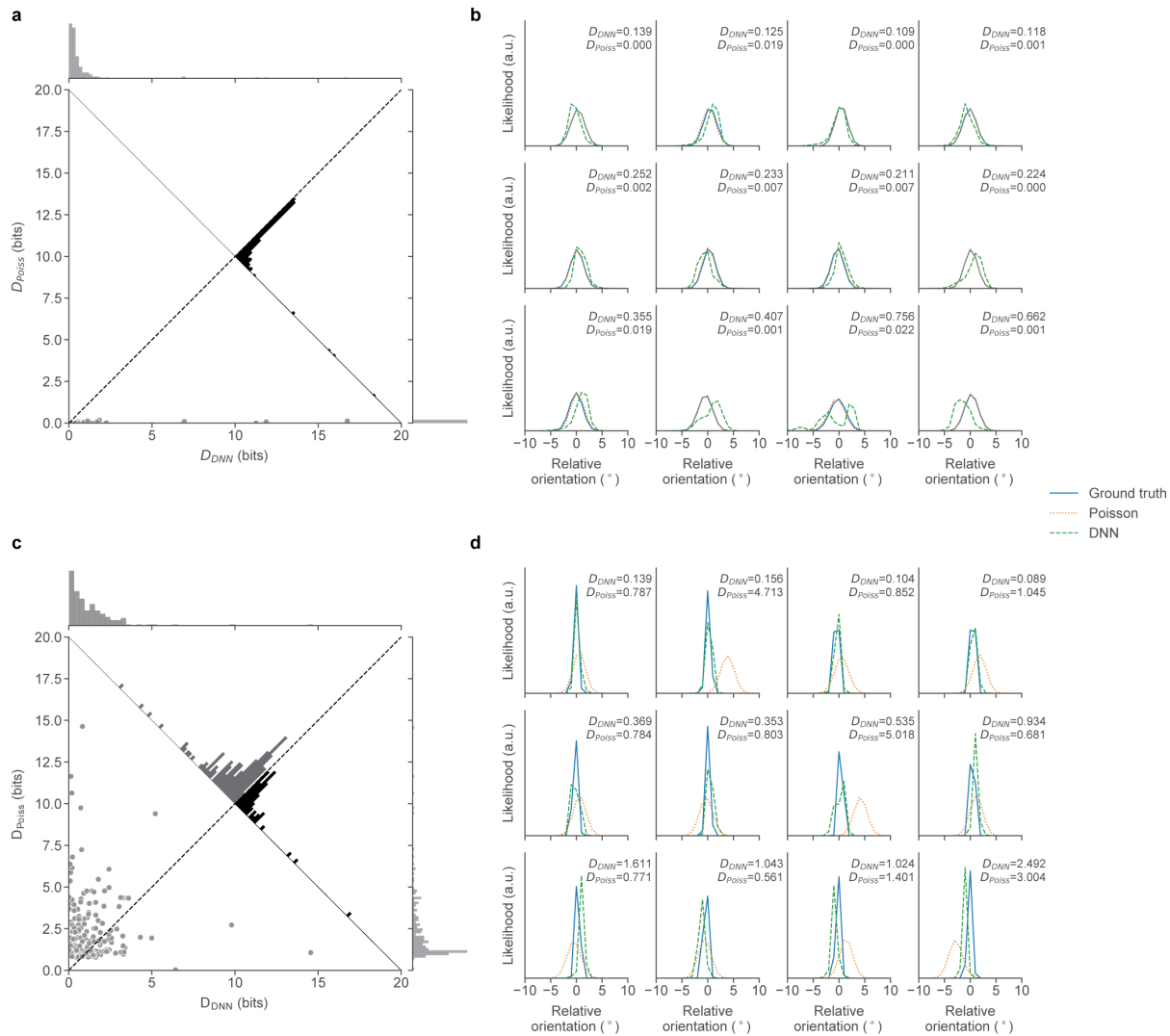222 iorally relevant stimulus uncertainty information, even when conditioned on the stimulus.

**Discussion**

224 Given the stochastic nature of the brain, repeated presentations of even identical stimuli elicit

225 variable responses. This noise creates an implicit code where multiple hypothetical stimuli are

226 consistent with a given response but with differing likelihoods. Here, we decoded trial-to-trial

227 likelihood functions from the population activity in visual cortex, and found that a model utilizing

228 the full likelihood function predicted the monkeys' choices better than alternative models that

229 ignore variations in the shape of the likelihood function. Our results provide the first population-

230 level evidence in support of the theoretical framework of probabilistic population coding, i.e. that

231 the brain performs Bayesian inference under a generative model of the neural activity.
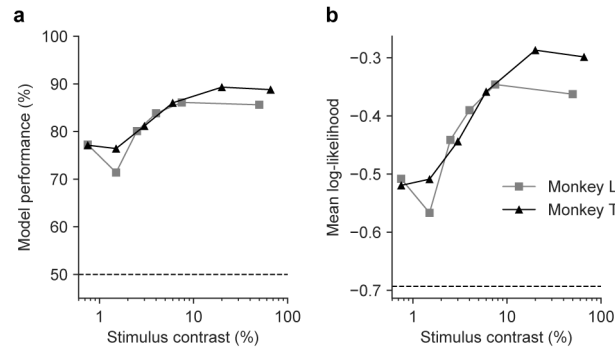
232 Our findings were made possible by recording from a large population simultaneously and by

233 using a task in which uncertainty is relevant to the animal. In addition, we developed a novel

234 method to decode the likelihood functions, which is based on a deep neural network and does not

235 rely on the strong parametric assumptions that have dominated previous work. Importantly, we

236 conditioned our analyses on the stimulus to rule out a confounding effect of the stimulus on the

237 observed relationship between the decoded likelihood function and the subject's decision. This

238 is critical because previous behavioral studies on cue combination and Bayesian integration, for

17

239 instance, always relied on varying stimulus features (e.g. contrast, blur, motion coherence) to ma-

240 nipulate uncertainty[7,8,16,34]. As a result, these studies cannot rule out that correlation between a

241 proposed encoding of uncertainty and the behavior may be confounded by the stimulus (Fig. 1b,c),

242 and therefore fail to provide rigorous assessment on the representation of uncertainty. Carefully

243 controlling for the effect of stimulus fluctuations allowed us to present rigorous evidence that the

244 trial-by-trial fluctuations in the likelihood functions carry behaviorally relevant stimulus uncer-
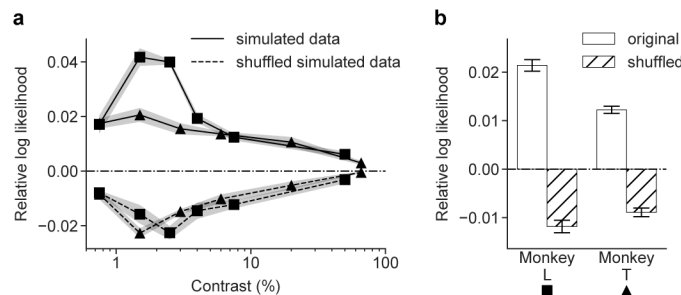
245 tainty information.

246 While the sensory likelihood function is a crucial building block for probabilistic computation in

247 the brain, fundamental questions remain regarding the nature of such computation. First, how do

248 downstream areas process the information contained in sensory likelihood functions to make better

249 decisions? Previous work has manually constructed neural networks for downstream computation

250 that relied heavily on the assumption of Poisson-like variability[9,10,17,35–37]. However, more recent

251 work has demonstrated that training generic shallow networks accomplishes the same goal without

252 the need for task-specific manual construction[38]. Second, does each area in a feedforward chain of

253 computation encode a likelihood function over its own variable, with the computation propagating

254 the uncertainty information from one variable to the next? For example, in our task, it is conceiv-

255 able that prefrontal cortex encodes a likelihood function over class that is derived from a likelihood

256 function over orientation coming in from V1. While answering these questions will require major

257 efforts, we expect that our findings regarding the neural basis of sensory uncertainty will help put

258 those efforts on a more solid footing.

**Extended Data Figure 1:** Performance of the likelihood functions decoded by DNN-based decoders. **a-b**, Results on independent Poisson population responses. **a**, KL divergence between the ground truth likelihood function and likelihood function decoded with: a trained DNN $D_{\text{DNN}}$ vs. independent Poisson distribution assumption $D_{\text{Poiss}}$. Each point is a single trial in the test set. The distributions of $D_{\text{DNN}}$ and $D_{\text{Poiss}}$ are shown at the top and right margins, respectively. The distribution of pair-wise difference between $D_{\text{DNN}}$ and $D_{\text{Poiss}}$ is shown on the diagonal. **b**, Example likelihood functions. The ground truth (solid blue), independent-Poisson based (dotted orange), and DNN-based (dashed green) likelihood functions are shown for selected trials from the test set. Four random samples (columns) were drawn from the top, middle and bottom 1/3 of trials sorted by the $D_{\text{DNN}}$ (rows). **c-d**, Same as in **a-b** but for simulated population responses with correlated Gaussian distribution where variance is scaled by the mean.

19

**Extended Data Figure 2:** Full-Likelihood Model performance. **a**, Model performance measured in proportions of trials correctly predicted by the model as a function of contrast for two monkeys (squares and triangles marks for Monkey L and T, respectively). On each trial, the class decision that would maximize the posterior $p(\hat{C}|\mathbf{r})$ was chosen to yield a concrete classification prediction. **b**, Same as in **a** but with performance measured as the trial-averaged log likelihood of the model. For **a** and **b**, dashed lines indicate the performance at chance ($50\%$ and $\ln(0.5)$, respectively).



**Extended Data Figure 3:** Expected model performance on simulated data using the trained Full-Likelihood Model as the ground truth. **a**, Average trial-by-trial performance of the Full-Likelihood Model relative to the Fixed-Uncertainty Model across contrasts on the simulated data, measured as the trial-averaged difference in the log likelihood. The results for the unshuffled and the shuffled simulated data are shown in solid and dashed lines, respectively. The squares and triangles mark Monkey L and T, respectively. **b**, Relative model performance summarized across all contrasts. Results are shown for each monkey and for unshuffled (solid) and shuffled (hatched) simulated data. For **a** and **b**, all data points are the means and error bar/shaded area indicate the standard deviation across the 5 simulation repetitions.

## Methods

**Experimental model and subject details** All behavioral and electrophysiological data were obtained from two healthy, male rhesus macaque (*Macaca mulatta*) monkeys (L and T) aged 10 and 7 years and weighting 9.5 and 15.1 kg, respectively. All experimental procedures complied with guidelines of the NIH and were approved by the Baylor College of Medicine Institutional Animal Care and Use Committee (permit number: AN-4367). Animals were housed individually in a room located adjacent to the training facility on a 12h light/dark cycle, along with around ten other monkeys permitting rich visual, olfactory, and auditory social interactions. Regular veterinary care and monitoring, balanced nutrition and environmental enrichment were provided by the Center for Comparative Medicine of Baylor College of Medicine. Surgical procedures on monkeys were conducted under general anesthesia following standard aseptic techniques.

**Stimulus presentation** Each visual stimulus was a single drifting oriented sinusoidal grating (spatial frequency: 2.79 cycles/degree visual angle, drifting speed: 0.028 cycles/s) presented through a circular aperture situated at the center of the screen. The size of the aperture was adjusted to cover receptive fields of the recorded populations, extending 2.14° and 2.86° of visual angle for Monkey L and Monkey T, respectively. The orientation and contrast of the stimulus were adjusted on a trial-by-trial basis as will be described later. The stimulus was presented on a CRT monitor (at a distance of 100 cm; resolution: 1600 × 1200 pixels; refresh rate: 100 Hz) using Psychophysics Toolbox[39]. The monitor was gamma-corrected to have a linear luminance response profile. Video cameras (DALSA genie HM640; frame rate 200Hz) with custom video eye tracking software developed in LabVIEW were used to monitor eye movements.

**Behavioral paradigm** On a given trial, the monkey viewed a drifting oriented grating with orientation $\theta$, drawn from one of two classes, each defined by a Gaussian probability distribution. Both distributions have a mean of $0°$ (grating drifting horizontally rightward, positive orientation corresponding to counter-clockwise rotation), but their standard deviations differed: $\sigma_1 = 3°$ for class 1 ($C = 1$) and $\sigma_2 = 15°$ for class 2 ($C = 2$). On each trial, the class was chosen randomly with equal probability, with the orientation of the stimulus then drawn from the corresponding distribution, $p(\theta|C)$. At the beginning of each recording session, at least three distinct values of contrasts were selected, and one of these values was chosen at random on each trial. Unlike more typical two-category tasks using distributions with identical variances but different means, optimal decision-making in our task requires the use of sensory uncertainty on a trial-by-trial basis[17].

Each trial proceeded as follows. A trial was initiated by a beeping sound and the appearance of a fixation target ($0.15°$ visual angle) in the center of the screen. The monkey fixated on the fixation target for 300 ms within $0.5°$–$1°$ visual angle. The stimulus then appeared at the center of the screen. After 500 ms, two colored targets (red and green) appeared to the left and the right of the grating stimulus (horizontal offset of $4.29°$ from the center with the target diameter of $0.71°$ visual angle), at which point the monkey saccaded to one of the targets to indicate their choice of class. For Monkey L, the grating stimulus was removed from the screen when the saccade target appeared, while for Monkey T, the grating stimulus remained on the screen until the subject completed the task by saccading to the target. The left-right configuration of the colored targets were varied randomly for each trial. Through training, the monkey learned to associate the red and the green targets with the narrow ($C = 1$) and the wide ($C = 2$) class distributions, respectively.

22

301  For illustrative clarity, we used blue to indicate $C = 2$ throughout this document. The monkey

302  received a juice reward for each correct response (0.10–0.15 mL).

303  During the training, the monkeys were first trained to perform the colored version of the task,

304  where the grating stimulus was colored to match the correct class $C$ for that trial (red for $C = 1$

305  and green for $C = 2$). Under this arrangement, the monkey simply learned to saccade to the target

306  matching the color of the grating stimulus, although the grating stimulus orientation information

307  was always present. As the training proceeded, we gradually removed the color from the stimulus,

308  encouraging the monkey to make use of the orientation information in the stimulus to perform the

309  task. Eventually, the color was completely removed, and at that point the monkey was performing

310  the full version of the task.

311  **Surgical Methods** Our surgical procedures followed a previously established approach[29,40,41].

312  Briefly, a custom-built titanium cranial headpost was first implanted for head stabilization under

313  general anesthesia using aseptic conditions in a dedicated operating room. After premedication

314  with Dexamethasone (0.25-0.5 mg/kg; 48 h, 24 h and on the day of the procedure) and atropine

315  (0.05 mg/kg prior to sedation), animals were sedated with a mixture of ketamine (10 mg/kg) and

316  xylazine (0.5 mg/kg). During the surgery, anesthesia was maintained using isoflurane (0.5-2%).

317  After the monkey was fully trained, we implanted a 96-electrode microelectrode array (Utah array,

318  Blackrock Microsystems, Salt Lake City, UT, USA) with a shaft length of 1 mm over parafoveal

319  area V1 on the right hemisphere. This surgery was performed under identical conditions as de-

320  scribed for headpost implantation. To ameliorate pain, analgesics were given for 7 days following

321  a surgery.

**Electrophysiological recording and data processing** The neural signals were pre-amplified at the head stage by unity gain preamplifiers (HS-27, Neuralynx, Bozeman MT, USA). These signals were then digitized by 24-bit analog data acquisition cards with 30 dB onboard gain (PXI-4498, National Instruments, Austin, TX) and sampled at 32 kHz. Broadband signals (0.5 Hz to 16 kHz) were continuously recorded using custom-built LabVIEW software for the duration of the experiment. Eye positions were tracked at 200 Hz using video cameras (DALSA genie HM640) with custom video eye tracking software developed in LabVIEW. The spike detection was performed offline according to a previously described method[27, 29, 40]. Briefly, a spike was detected when the signal on a given electrode crossed a threshold of five times the standard deviation of the corresponding electrode. To avoid artificial inflation of the threshold in the presence of a large number of high-amplitude spikes, we used a robust estimator of the standard deviation[42], given by $\text{median}(|x|)/0.6745$. Spikes were aligned to the center of mass of the continuous waveform segment above half the peak amplitude. Code for spike detection is available online at `https://github.com/atlab/spikedetection`. In this study, the term "multiunit" refers to the set of all spikes detected from a single channel (i.e. electrode) of the Utah array, and all analyses in the main text were performed on multiunits. For each multiunit, the total number of spikes during the 500 ms of pre-target stimulus presentation, $r_i$ for the $i^{\text{th}}$ unit, was used as the measure of the multiunit's response for a single trial. The population response $\mathbf{r}$ is the vector of spike counts for all 96 multiunits.

**Dataset and inclusion criteria.** We recorded a total of 61 and 71 sessions from Monkey L and T, for a total of 112,072 and 193,629 trials, respectively. We removed any trials with electrophys-

24

343  iology recordings contaminated by noise in the recording devices (e.g. poor grounding connector

344  resulting in movement noise) or equipment failures. To do so, we established the following trial

345  inclusion criteria:

1. The total spike counts $r_t = \sum_i r_i$ across all channels should fall within the $\pm 4\sigma_{\mathrm{adj}}$ from the median total spike counts across all trials from a single session. $\sigma_{\mathrm{adj}}$ is the standard deviation of the total spike count distribution robustly approximated using the interquartile range IQR as follows: $\sigma_{\mathrm{adj}} = \frac{\mathrm{IQR}}{1.35}$.

2. For at least 50% of all units, the observed $i^{\mathrm{th}}$ unit spike count $r_i$ for the trial should fall within a range defined as: $|r_i - \mathrm{MED}_i| \leq 1.5 \cdot \mathrm{IQR}_i$, where $\mathrm{MED}_i$ and $\mathrm{IQR}_i$ are the median and interquartile ranges of the $i^{\mathrm{th}}$ unit spike counts distribution throughout the session, respectively.

354  We only included trials that satisfied both of the above criteria in our analysis. Empirically, we

355  found the above criteria to be effective in catching obvious anomalies in the spike data while

356  introducing minimal bias into the data. After the application of the criteria, we were left with

357  110,695 and 192,631 trials for Monkey L and T, thus retaining 98.77% and 99.48% of the total

358  trials, respectively. While this selection criteria allowed us to remove apparent anomaly in the

359  data, we found that the main findings described in this paper were not sensitive to the precise

360  definition of the inclusion criteria.

361  During each recording session, stimuli were presented under three or more contrast values. In all

362  analyses to follow, we studied the trials from distinct contrast separately for each recording session,

25

363 and we refer to this grouping as a "contrast-session".

**Receptive field mapping** On the first recording session for each monkey, the receptive field was mapped using spike-triggered averaging of the multiunit responses to a white noise random dot stimulus. The white noise stimulus consisted of square dots of size $0.29°$ of visual angle presented on a uniform gray background, with randomly varying location and color (black or white) every 30 ms for 1 second. We adjusted the size of the grating stimulus as necessary to ensure that the stimulus covers the population receptive field entirely.

**Likelihood decoding** Given the population activity $\mathbf{r}$ in response to an orientation $\theta$, we aimed to decode uncertainty information in the form of a likelihood function $\mathcal{L}(\theta) \equiv p(\mathbf{r}|\theta)$, as a function of $\theta$. This may be computed through the knowledge of the generative relation leading from $\theta$ to $\mathbf{r}$—that is, by describing the underlying orientation conditioned probability distribution over $\mathbf{r}$, $p(\mathbf{r}|\theta)$. This procedure is typically approximated by making rather strong assumptions about the form of the density function, for example by assuming that neurons fire independently and each neuron fires according to the Poisson distribution[21]. Under this approach, the expected firing rates (i.e. tuning curves) of the $i^{\text{th}}$ neuron $\mathbb{E}[r_i|\theta] = f_i(\theta)$ must be approximated as well, for example by fitting a parametric function (e.g. von Mises tuning curves[43]) or employing kernel regression[21]. While these approaches have proven useful, the effect of the strong and likely inaccurate assumptions on the decoded likelihood function remains unclear. Ideally, we can more directly estimate the likelihood function $\mathcal{L}(\theta)$ without having to make strong assumptions about the underlying conditional probability distribution over $\mathbf{r}$.

26

383 To this end, we employed a deep neural network (DNN)[18] to directly approximate the likelihood

384 function over the stimulus orientation, $\theta$, from the recorded population response $\mathbf{r}$. Here we present

385 a brief derivation that serves as the basis of the network design and training objective. Let us

386 assume that $m$ multiunits were recorded simultaneously in a single recording session, so that $\mathbf{r} \in$

387 $\mathbb{R}^m$. To make the problem tractable, we bin the stimulus orientation $\theta$ into $n$ distinct values, $\theta_1$

388 to $\theta_n$ (the derivation holds in general for arbitrarily fine binning of the orientation). With this, the

389 likelihood function can be captured by a vector $\mathbf{L} \in \mathbb{R}^n$ where $\mathbf{L}_i = \mathcal{L}(\theta_i)$. Let us assume that

390 we can train some DNN to learn a mapping $f$ from the population response $\mathbf{r}$ to the log of the

391 likelihood function $\mathbf{L}$ up to a constant offset $b$. That is, $f : \mathbb{R}^m \mapsto \mathbb{R}^n$,

$$\mathbf{r} \mapsto f(\mathbf{r}) = \log \mathbf{L} + b(\mathbf{r}) = \log p(\mathbf{r}|\theta) + b(\mathbf{r}) \tag{1}$$

392 for some scalar function $b \in \mathbb{R}$. As the experimenter, we know the distribution of the stimulus

393 orientation, $\mathbf{p}_\theta \in \mathbb{R}^n$, where $\mathbf{p}_{\theta,i} = p(\theta_i)$. We combine $f(\mathbf{r})$ and $\mathbf{p}_\theta$ to compute the log posterior

394 over stimulus orientation $\theta$ up to some scalar value $b'(\mathbf{r})$,

$$\mathbf{z}(\mathbf{r}) \equiv \log \mathbf{p}_\theta + f(\mathbf{r}) = \log p(\theta|\mathbf{r}) + b'(\mathbf{r}) \tag{2}$$

395 We finally take the softmax of $\mathbf{z}(\mathbf{r})$, and recover the normalized posterior function $\mathbf{q}(\mathbf{r}) \equiv \text{softmax}(\mathbf{z}(\mathbf{r}))$

396 where,

27

$$\mathbf{q}_i(\mathbf{r}) = \frac{e^{\mathbf{z}_i(\mathbf{r})}}{\sum_j e^{\mathbf{z}_j(\mathbf{r})}} \tag{3}$$

$$= \frac{e^{b'(\mathbf{r})} p(\theta = \theta_i | \mathbf{r})}{e^{b'(\mathbf{r})} \sum_j p(\theta = \theta_j | \mathbf{r})} \tag{4}$$

$$= p(\theta = \theta_i | \mathbf{r}) \tag{5}$$

397 Overall, $\mathbf{q}(\mathbf{r}) = \text{softmax}(\log \mathbf{p} + f(\mathbf{r}))$.

398 The goal then is to train the DNN $f(\mathbf{r})$ such that the overall function $\mathbf{q}(\mathbf{r})$ matches the posterior

399 over the stimulus, $\mathbf{p}(\mathbf{r})$ where $\mathbf{p}_i(\mathbf{r}) = p(\theta = \theta_i | \mathbf{r})$ based on the available data. This in turn

400 allows the network output $f(\mathbf{r})$ to approach the log of the likelihood function $\mathbf{L}$, up to a constant

401 $b(\mathbf{r})$. For 1-out-of-$n$ classification problems, minimizing the cross-entropy between $\mathbf{q}(\mathbf{r})$ and the

402 stimulus orientation $\theta$ for a given $\mathbf{r}$ lets the overall function $\mathbf{q}(\mathbf{r})$ approach the true posterior $\mathbf{p}(\mathbf{r})$,

403 as desired[44,45]. To show this, let us start by minimizing the difference between the model estimated

404 posterior $\mathbf{q}(\mathbf{r})$ and the true posterior $\mathbf{p}(\mathbf{r})$ over the distribution of $\mathbf{r}$. We do this by minimizing

405 the loss $L$ defined as the expected value of the Kullback-Leibler divergence[46] between the two

406 posteriors:

28

$$L(W) = \mathbb{E}_{\mathbf{r}}\left[D_{KL}(\mathbf{p}||\mathbf{q})\right] \tag{6}$$

$$= \mathbb{E}_{\mathbf{r}}\left[\mathbb{E}_{\theta|\mathbf{r}}\left[\log\frac{p(\theta|\mathbf{r})}{q(\theta|\mathbf{r},W)}\right]\right] \tag{7}$$

$$= \mathbb{E}_{\mathbf{r},\theta}\left[\log\frac{p(\theta|\mathbf{r})}{q(\theta|\mathbf{r},W)}\right] \tag{8}$$

$$= -\mathbb{E}_{\mathbf{r},\theta}\left[\log q(\theta|\mathbf{r},W)\right] - H(\theta|\mathbf{r}) \tag{9}$$

where $p(\theta = \theta_i|\mathbf{r}) \equiv \mathbf{p}_i(\mathbf{r})$, $q(\theta = \theta_i|\mathbf{r},W) \equiv \mathbf{q}_i(\mathbf{r},W)$, $W$ is a collection of all trainable parameters in the network, and $H(\theta|\mathbf{r})$ is the conditional entropy of $\theta$ conditioned on $\mathbf{r}$, which is an unknown but a fixed quantity with respect to $W$ and the data distribution. Here we used the notation $\mathbf{q}(\mathbf{r},W)$ to highlight the dependence of the network estimated posterior $\mathbf{q}(\mathbf{r})$ on the network parameters $W$. We can redefine the loss, $L^*$, only leaving the terms that depends on the trainable parameters $W$, and then apply a Monte Carlo method[47] to approximate the loss from samples:

$$L^*(W) = -\mathbb{E}_{\mathbf{r},\theta}\left[\log q(\theta|\mathbf{r},W)\right] \tag{10}$$

$$\approx -\frac{1}{N}\sum_i \log q(\theta^{(i)}|\mathbf{r}^{(i)},W) \tag{11}$$

where $\left(\theta^{(i)}, \mathbf{r}^{(i)}\right)$ are samples drawn from a training set for the network. Eq. 11 is precisely the definition of the cross-entropy as we set out to show.

Therefore, by optimizing the overall function $\mathbf{q}(\mathbf{r})$ to match the posterior distribution through

417   the use of cross-entropy loss, the network output $f(\mathbf{r})$ can approximate the log of the likelihood

418   function $\mathcal{L}(\theta)$ for each $\mathbf{r}$ up to an unknown constant $b(\mathbf{r})$. Because we do not know the value of

419   $b(\mathbf{r})$, the network will not learn to recover the underlying generative function linking from $\theta$ to $\mathbf{r}$,

420   $p(\mathbf{r}|\theta)$.

421   As an example, consider a neural population with responses that follows a Poisson-like distri-

422   bution (i.e. a version of the exponential distribution with linear sufficient statistics[9, 10]). Learn-

423   ing a decoder for such population responses occurs as a special case of training a DNN-based

424   likelihood decoder. For Poisson-like variability, the stimulus-conditioned distribution over $\mathbf{r}$ is

425   $p(\mathbf{r}|\theta) = \phi(\mathbf{r})e^{\mathbf{h}^\top(\theta)\mathbf{r}}$. The log likelihood function is then $\log \mathbf{L} = \log \phi(\mathbf{r}) + \mathbf{H}^\top \mathbf{r}$, where $\mathbf{H}$ is a

426   matrix whose $i^{\text{th}}$ column is $\mathbf{h}(\theta_i)$. If we let $f(\mathbf{r}) = \mathbf{H}^\top \mathbf{r}$, then $f(\mathbf{r}) = \log \mathbf{L} + b(\mathbf{r})$ as desired, for

427   $b(\mathbf{r}) = -\log \phi(\mathbf{r})$. Hence, if we used a simple fully connected network, training the network is

428   equivalent to fitting the kernel function $\mathbf{h}(\theta)$ of the Poisson-like distribution.

429   In this work, we modeled the mapping $f(\mathbf{r})$ as a DNN with two hidden layers[19], consisting of

430   two repeating blocks of a fully connected layer followed by a rectified linear unit (ReLU)[18] and

431   a drop-out layer[48], and a fully connected readout layer with no output nonlinearity (Figure 3c).

432   To encourage smoother likelihood functions, we added an $L_2$ regularizer on $\log \mathbf{L}$ filtered with

433   a Laplacian filter of the form $\mathbf{h} = [-0.25, 0.5, -0.25]$. Therefore, the training loss included the

434   term:

$$R = \gamma \sum_i \mathbf{u}_i^2 \tag{12}$$

435 for $\mathbf{u} = (\log \mathbf{L}) * \mathbf{h}$, where $*$ denotes convolution operation, $\mathbf{u}_i$ is the $i^{\text{th}}$ element of the filtered log

436 likelihood function $\mathbf{u}$, and $\gamma$ is the weight on the smoothness regularizer.

437 We trained a separate instance of the network for each contrast-session. During the training, each

438 contrast-session was randomly split in proportions of 80% / 20% to yield the training set and

439 the validation set, respectively. The stimulus orientation $\theta$ was binned into integers in the range

440 $[-45°, 45°]$, and we excluded trials with orientations outside this range. This led to the exclusion of

441 157 out of 110,695 trials (0.14%) and 254 out of 192,631 trials (0.13%) for Monkey L and T data,

442 respectively. The network was trained on the training set, and its performance on the validation

443 set was monitored to perform early stopping[49], and subsequently hyperparameter selection. For

444 early stopping, we computed the mean squared error (MSE) between the maximum-a-posteriori

445 (MAP) readout of the network output posterior $\mathbf{q}$ and the stimulus orientation $\theta$ on the validation

446 set, and the training was terminated (early-stopped) if MSE failed to improve over 400 consecutive

447 epochs, where each epoch is defined as one full pass through the training set. Upon early stopping,

448 the parameter set that yielded the best validation set MSE during the course of the training was

449 restored. Once the training was complete, the trained network was evaluated on the validation

450 set to yield the final score, which served as the basis for our hyperparameter selections. The

451 values of hyperparameters for the networks, including the size of hidden layers, the weight on the

452 likelihood function smoothness regularizer $\gamma$, and drop-out rates during the training were selected

453 by performing a grid search over candidate values to find the combination that yielded the best

454 validation set score for each contrast-session instance of the network.

**Decoding likelihood functions from known response distributions** To assess the effectiveness of the DNN-based likelihood decoding method described above, we simulated neural population responses with known noise distributions, trained DNN decoders on the simulated population responses, and compared the decoded likelihood functions to the ground-truth likelihood functions obtained by inverting the known generative model for the responses. We also compared the quality of the DNN-decoded likelihood functions to those decoded by assuming independent Poisson distribution on the population responses, as done in previous work[14, 16, 20, 21, 23].

We simulated the activities of a population of 96 multiunits $\mathbf{r}_{\mathrm{sim}}$ responding to the stimulus orientation $\theta$ drawn from the the distribution defined for our task such that:

$$p(\theta) = \frac{1}{2}\mathcal{N}(\theta; 0, \sigma_1^2) + \frac{1}{2}\mathcal{N}(\theta; 0, \sigma_2^2) \tag{13}$$

where $\sigma_1 = 3°$ and $\sigma_2 = 15°$.

We modeled the expected response of $i^{\mathrm{th}}$ unit to $\theta$—that is, the tuning function $f_i(\theta)$—with a Gaussian function:

$$f_i(\theta) = Ae^{-\frac{(\theta - \mu_{\mathrm{sim},i})^2}{2\sigma_{\mathrm{sim}}^2}} \tag{14}$$

For the simulation, we have set $A = 6$ and $\sigma_{\mathrm{sim}} = 21°$. We let the mean of the Gaussian tuning curves for the 96 units to uniformly tile the stimulus orientation between $-40°$ and $40°$. In other words,

32

$$\mu_{\text{sim},i} = -40 + \frac{16}{19}(i-1) \tag{15}$$

for $i \in [1, 96]$.

For any given trial with a drawn orientation $\theta$, the population response $\mathbf{r}_{\text{sim}}$ was then generated under two distinct models of distributions. In the first case, the population responses were drawn from an independent Poisson distribution as is commonly assumed in many works:

$$p(\mathbf{r}_{\text{sim}}|\theta) = \prod_i \text{Poiss}(r_{\text{sim},i}; f_i(\theta)) \tag{16}$$

$$= \prod_i \frac{f_i(\theta)^{r_{\text{sim},i}} e^{-f_i(\theta)}}{r_{\text{sim},i}!} \tag{17}$$

In the second case, the population responses were drawn from a multivariate Gaussian distribution with covariance matrix $\Sigma \in \mathbb{R}^{96 \times 96}$ that scales with the mean response of the population. That is:

$$p(\mathbf{r}_{\text{sim}}|\theta) = \mathcal{N}(\mathbf{r}_{\text{sim}}; \mathbf{f}(\theta), \Sigma(\theta)) \tag{18}$$

for

$$\Sigma(\theta) = (\mathbf{f}^{1/2}(\theta))^\top C(\mathbf{f}^{1/2}(\theta)) \tag{19}$$

477 where $\mathbf{f}^{1/2}(\theta) \in \mathbb{R}^{96}$ such that $\mathbf{f}_i^{1/2}(\theta) = \sqrt{f_i(\theta)}$, and $C \in \mathbb{R}^{96 \times 96}$ is a correlation matrix. Under

478 this distribution, the variance of any unit's response scales linearly with its mean just as in the case

479 of the Poisson distribution, but the population responses can be highly correlated depending on the

480 choice of the correlation matrix $C$. For the simulation, we randomly generated a correlation matrix

481 with the average units correlation of 0.227.

482 For each case of the distribution, we simulated population responses for the total of 1200 trials.

483 Among these, 200 trials were set aside as the test set. We trained the DNN-based likelihood

484 decoder on the remaining 1000 trials, splitting them further into 800 and 200 trials as the training

485 and validation set, respectively. We followed the exact DNN training and hyperparameter selection

486 procedure as described earlier.

487 For comparison, we also decoded the likelihood function from the population response $\mathbf{r}_{\text{sim}}$ under

488 the assumption of independent Poisson variability, regardless of the "true" distribution. Each unit's

489 responses over the 1000 trials were fitted separately with a Gaussian tuning curve (Eq. 14). The

490 parameters of the tuning curve $A_i$, $\mu_i$ and $\sigma_{\text{sim, i}}$ were obtained by minimizing the least square

491 difference between the Gaussian tuning curve and the observed $i^{\text{th}}$ unit's responses $(\theta, r_{\text{sim},i})$ using

492 `least_squares` function from Python SciPy optimization library.

493 The ground-truth likelihood function $p(\mathbf{r}_{\text{sim}}|\theta)$ was computed for each simulated trial according

494 to the definition of the distribution as found in Eq. 16 for the independent Poisson population or

495 Eq. 18 for the mean scaled correlated Gaussian population.

496 We then assessed the quality of the decoded likelihood functions under the independent Pois-

$_{497}$ son model $\mathcal{L}_{\mathrm{Poiss}}(\theta)$ and under the DNN model $\mathbf{L}_{\mathrm{DNN}}$ by computing their Kullback-Leibler (KL)

$_{498}$ divergence[46] from the ground-truth likelihood function $\mathcal{L}_{\mathrm{gt}}(\theta)$, giving rise to $D_{\mathrm{Poiss}}$ and $D_{\mathrm{DNN}}$, re-

$_{499}$ spectively. All continuous likelihood functions ($\mathcal{L}_{\mathrm{gt}}$ and $\mathcal{L}_{\mathrm{Poiss}}$) were sampled at orientation $\theta$ where

$_{500}$ $\theta \in \mathbb{Z}$ and $\theta \in [-45°, 45°]$, giving rise to discretized likelihood functions $\mathbf{L}_{\mathrm{gt}}$ and $\mathbf{L}_{\mathrm{Poiss}}$ matching

$_{501}$ the dimensionality of the discretized likelihood function $\mathbf{L}_{\mathrm{DNN}}$ computed by the DNN. We then

$_{502}$ computed the KL divergence as:

$$D_{\mathrm{Poiss}} = \sum_i \log \frac{\mathbf{L}_{\mathrm{gt},i}}{\mathbf{L}_{\mathrm{Poiss},i}} \mathbf{L}_{\mathrm{gt},i} \tag{20}$$

$_{503}$ and

$$D_{\mathrm{DNN}} = \sum_i \log \frac{\mathbf{L}_{\mathrm{gt},i}}{\mathbf{L}_{\mathrm{DNN},i}} \mathbf{L}_{\mathrm{gt},i} \tag{21}$$

$_{504}$ We computed the KL divergence for both models across all 200 trials in the test set for both simu-

$_{505}$ lated population distributions (Extended Data Fig. 1a,c). When the simulated population distribu-

$_{506}$ tion was independent Poisson, then $D_{\mathrm{Poiss}} < D_{\mathrm{DNN}}$ for all test set trials (Extended Data Fig. 1a),

$_{507}$ indicating that $\mathbf{L}_{\mathrm{Poiss}}$ better approximated $\mathbf{L}_{\mathrm{gt}}$ overall than $\mathbf{L}_{\mathrm{DNN}}$. However, $\mathbf{L}_{\mathrm{DNN}}$ still closely ap-

$_{508}$ proximated $\mathbf{L}_{\mathrm{gt}}$ (Extended Data Fig. 1b).

$_{509}$ When the simulated population distribution was mean scaled correlated Gaussian, $\mathbf{L}_{\mathrm{DNN}}$ better

$_{510}$ approximated $\mathbf{L}_{\mathrm{gt}}$ than $\mathbf{L}_{\mathrm{Poiss}}$ on the majority of the trials (Extended Data Fig. 1c). Furthermore,

$_{511}$ $\mathbf{L}_{\mathrm{Poiss}}$ provided qualitatively worse fit to the $\mathbf{L}_{\mathrm{gt}}$ for the simulated correlated Gaussian distribution

35

512 compared to the fit of $\mathbf{L}_{\mathrm{DNN}}$ to $\mathbf{L}_{\mathrm{gt}}$ for the simulated independent Poisson distribution (Extended

513 Data Fig. 1b,d).

514 Overall, the simulation results suggest that (1) when the form of the underlying population dis-

515 tribution is known (such as in the case of independent Poisson population), more accurate like-

516 lihood functions can be decoded by directly using the knowledge of the population distribution

517 than through the DNN-based likelihood decoder, but (2) when the form of the underlying distri-

518 bution is unknown (such as in the case of the mean scaled correlated Gaussian distribution), then

519 a DNN-based likelihood decoder can yield much more accurate likelihood functions than if one

520 was to employ a wrong assumption about the underlying distribution in decoding likelihood func-

521 tions, and (3) a DNN-based likelihood decoder can provide reasonable estimate of the likelihood

522 function across wide range of underlying distributions. Because the true underlying population

523 distribution is hardly ever known to the experimenter, we believe that our DNN-based likelihood

524 decoder stands as the most flexible method in decoding likelihood functions from the population

525 responses to stimuli.

526 **Mean and standard deviation of likelihood function** For uses in the subsequent analyses, we

527 computed the mean and the standard deviation of the likelihood function by treating the likelihood

528 function as an unnormalized probability distribution:

$$\mu_L = \frac{\int \theta \mathcal{L}(\theta)\, d\theta}{\int \mathcal{L}(\theta)\, d\theta} \tag{22}$$

36

$$\sigma_L = \sqrt{\frac{\int (\theta - \mu_L)^2 \mathcal{L}(\theta) \, d\theta}{\int \mathcal{L}(\theta) \, d\theta}} \tag{23}$$

529 We took the $\mu_L$ and $\sigma_L$ to be the point estimate of the stimulus orientation and the measure of

530 the spread of the likelihood function, respectively, used in all subsequent analyses. Although not

531 presented here, we performed the following analyses with other point estimates of the stimulus

532 orientation such as the orientation at the maximum of the likelihood function and the median of

533 the likelihood functions, and observed that models with mean of the likelihood function as the

534 point estimate performed the best.

535 **Decision-making models** Given the hypothesized representation of the stimulus and its uncer-

536 tainty in the form of the likelihood function $\mathcal{L}(\theta) \equiv p(\mathbf{r}|\theta)$, the monkey's trial-by-trial decisions

537 were modeled based on the assumption that the monkey computes the posterior probability over

538 the two classes $C = 1$ and $C = 2$, and utilizes this information in making decisions—that is,

539 in accordance to a model of a Bayesian decision maker. The orientation distributions for the two

540 classes are $p(\theta|C = 1) = \mathcal{N}(\theta; \mu, \sigma_1^2)$ and $p(\theta|C = 2) = \mathcal{N}(\theta; \mu, \sigma_2^2)$ with $\mu = 0$ and $\sigma_1 = 3°$ and

541 $\sigma_2 = 15°$ where $\mathcal{N}(\theta; \mu, \sigma^2)$ denotes a Gaussian distribution over $\theta$ with mean $\mu$ and variance $\sigma^2$.

542 The posterior ratio $\rho$ for the two classes is:

37

$$\rho = \frac{p(C = 2|\mathbf{r})}{p(C = 1|\mathbf{r})} \tag{24}$$

$$= \frac{p(C = 2) \int p(\mathbf{r}|\theta)p(\theta|C = 2)\, d\theta}{p(C = 1) \int p(\mathbf{r}|\theta)p(\theta|C = 1)\, d\theta} \tag{25}$$

$$= \frac{p(C = 2) \int \mathcal{L}(\theta)\mathcal{N}(\theta; \mu, \sigma_2^2)\, d\theta}{p(C = 1) \int \mathcal{L}(\theta)\mathcal{N}(\theta; \mu, \sigma_1^2)\, d\theta} \tag{26}$$

A Bayes-optimal observer should select the class with the higher probability—a strategy known as maximum-a-posteriori (MAP) decision-making:

$$\hat{C} = \underset{C}{\operatorname{argmax}}\, p(C|\mathbf{r}) \tag{27}$$

where $\hat{C}$ is the subject's decision. However, according to the posterior probability matching strategy[50,51], the decision of subjects on certain tasks are better modeled as sampling from the posterior probability:

$$p(\hat{C}) = p(C = \hat{C}|\mathbf{r}) \tag{28}$$

To capture either decision-making strategy, we modeled the subject's classification decision probability ratio as follows:

$$\frac{p(\hat{C} = 2)}{p(\hat{C} = 1)} = \left(\frac{p(C = 2|\mathbf{r})}{p(C = 1|\mathbf{r})}\right)^\alpha = \rho^\alpha \tag{29}$$

38

$_{550}$ where $\alpha \in \mathbb{R}^+$. When $\alpha = 1$, the decision-making strategy corresponds to the posterior probability

$_{551}$ matching while $\alpha = \infty$ corresponds to the MAP strategy[51]. We fitted the value of $\alpha$ for each

$_{552}$ contrast-session during the model fitting to capture any variation of the strategy. Furthermore, we

$_{553}$ incorporated a lapse rate $\lambda$, a fraction of trials on which the subject does not pay attention and

$_{554}$ makes a random decision. Hence, the final probability that the subject selects the class $C = 1$ was

$_{555}$ modeled as:

$$p(\hat{C} = 1) = (1 - \lambda)\frac{1}{1 + \rho^\alpha} + 0.5\lambda \tag{30}$$

$$= (1 - \lambda)\left[1 + \left(\frac{p(C = 2)\int \mathcal{L}(\theta)\mathcal{N}(\theta; \mu, \sigma_2^2)\, d\theta}{p(C = 1)\int \mathcal{L}(\theta)\mathcal{N}(\theta; \mu, \sigma_1^2)\, d\theta}\right)^\alpha\right]^{-1} + 0.5\lambda \tag{31}$$

$$= (1 - \lambda)\left[1 + \left(\frac{(1 - p(C = 1))\int \mathcal{L}(\theta)\mathcal{N}(\theta; \mu, \sigma_2^2)\, d\theta}{p(C = 1)\int \mathcal{L}(\theta)\mathcal{N}(\theta; \mu, \sigma_1^2)\, d\theta}\right)^\alpha\right]^{-1} + 0.5\lambda \tag{32}$$

$_{556}$ For each contrast-session, we fitted the above Bayesian decision model to the monkey's decisions

$_{557}$ by fitting the four parameters: $\mu$, $p(C = 1)$, $\alpha$, and $\lambda$. Fitting $\mu$ (the center of stimulus orientation

$_{558}$ distributions) and $p(C = 1)$ (prior over class) allowed us to capture the bias in the stimulation

$_{559}$ distribution that the subject may have acquired errorneously during the training, and fitting $\alpha$ and

$_{560}$ $\lambda$ allowed for the model to match the decision-making strategy employed by the subject.

$_{561}$ Utilizing the likelihood function $\mathcal{L}(\theta)$ decoded from the V1 population response via the DNN in

$_{562}$ Eq. 32 gave rise to the Full-Likelihood Model that made use of all information including the trial-

$_{563}$ by-trial uncertainty information as captured by the shape of the likelihood function. Alternatively,

$_{564}$ we approximated the trial-by-trial decoded likelihood function for a contrast-session with a (fitted)

39

fixed width Gaussian function whose mean matched that of the normalized decoded likelihood function. This Fixed-Uncertainty Model effectively discarded all trial-by-trial fluctuations in the uncertainty as captured by the shape of the likelihood function but preserved the point estimate of the stimulus orientation $\hat{\theta}$ (i.e. mean of the likelihood function). That is:

$$\hat{\mathcal{L}}(\theta) = \mathcal{N}(\theta; \hat{\theta}, \sigma_c^2) \tag{33}$$

where $\hat{\theta} = \mu_L$ (Eq. 22). For each contrast-session, different values were fitted for the width of the Gaussian likelihood function, $\sigma_c$. Therefore, the Fixed-Uncertainty Model had 5 parameters to be fitted in total: $\mu$, $p(C = 1)$, $\alpha$, $\lambda$, and $\sigma_c$.

**Model fitting and model comparison** We used 10-fold cross-validation to fit and evaluate both decision models, separately for each contrast-session. We divided all trials from a given contrast-session randomly into 10 equally sized subsets, $B_1, B_2, \ldots, B_i, \ldots, B_{10}$ where $B_i$ is the $i^{\text{th}}$ subset. We then held out a single subset $B_i$ as the test set, and trained the decision-making model on the remaining 9 subsets combined together, serving as the training set. The predictions and the performance of the trained model on the held out test set $B_i$ was then reported. We repeated this 10 times, iterating through each subset as the test set, training on the remaining subsets.

The decision models were trained to minimize the negative log likelihood on the subject's decision across all trials in the training set:

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \left( -\log \prod_i p(\hat{C} = \hat{C}_i | M, \Theta) \right) \tag{34}$$

$$= \underset{\Theta}{\operatorname{argmin}} \left( -\sum_i \log p(\hat{C} = \hat{C}_i | M, \Theta) \right) \tag{35}$$

where $\Theta$ is the collection of the parameters for the decision-making model $M$ and $\hat{C}_i$ is the subject's decision on the $i^{\text{th}}$ trial in the training set. The term $p(\hat{C}|M, \Theta)$ is given by the Eq. 32 with either the unmodified $\mathcal{L}(\theta)$ in the Full-Likelihood Model or a Gaussian approximation to $\mathcal{L}(\theta)$ in the Fixed-Uncertainty Model.

The optimizations were performed using three algorithms: `fmincon` and `ga` from MATLAB's optimization toolbox and Bayesian Adaptive Direct Search (BADS)[52]. When applicable, the optimization was repeated with 50 or more random parameter initializations. For each cross-validation fold, we retained the parameter combination $\hat{\Theta}$ that yielded the best training set score (i.e. lowest negative log likelihood) among all optimization runs across different algorithms and parameter initializations. We subsequently tested the model $M$ with the best training set parameter $\hat{\Theta}$ and reported the score on the test set. For each contrast-session, all analyses on the trained model presented in the main text were performed on the aggregated test sets scores.

**Likelihood shuffling analysis** To assess the contribution of the trial-by-trial fluctuations in the decoded likelihood functions in predicting the animal's decisions under the Full-Likelihood Model, for each contrast-session we shuffled the likelihood functions among trials in the same stimulus orientation bin, while maintaining the trial to trial relationship between the point estimate of the stimulus orientation (i.e. mean of the normalized likelihood) and the perceptual decision. Specif-

41

598 ically, we binned trials to the nearest orientation degree such that each bin was centered at an

599 integer degree (i.e. bin center $\in \mathbb{Z}$) with the bin width of $1°$. We then shuffled the likelihood func-

600 tions among trials in the same orientation bin. This effectively removed the stimulus orientation

601 conditioned correlation between the likelihood function and the subject's classification $\hat{C}$, while

602 preserving the expected likelihood function for each stimulus orientation.

603 However, we were specifically interested in decoupling the uncertainty information contained in

604 the shape of the likelihood function from the decision while minimally disrupting the trial-by-trial

605 correlation between the point estimate of the stimulus orientation and the subject's classification

606 decision. To achieve this, for each trial, the newly assigned likelihood function was shifted such

607 that the mean of the normalized likelihood function, $\mu_L$ (Eq. 22), remained the same for each

608 trial before and after the likelihood shuffling (Fig. 5c). This allowed us to specifically assess the

609 effect of distorting the shape of the likelihood function conditioned on both the (binned) stimulus

610 orientation and the point estimate of the stimulus orientation (i.e. $\mu_L$) (Fig. 5c). We then trained

611 both the Full-Likelihood Model and the Fixed-Uncertainty Model on the shuffled data, following

612 the exact procedure used when training on the original (unshuffled) data.

613 **Classification simulation** We computed the expected effect size of the model fit difference be-

614 tween the Full-Likelihood Model and the Fixed-Uncertainty Model by generating simulated data

615 using the trained Full-Likelihood Model as the ground truth. Specifically, for each trial for each

616 contrast-session, we computed the probability of responding $\hat{C} = 1$ from Eq. 32, utilizing the

617 full decoded likelihood function $\mathcal{L}(\theta)$ for the given trial, and sampled a classification decision

618 from that probability. This procedure yielded simulated data where all monkeys' decisions were

42

619 replaced by decisions made by the trained Full-Likelihood Models. We repeated this procedure

620 5 times, thereby producing 5 sets of simulated data. For each set of simulated data, we trained

621 the two decision-making models (Full-Likelihood Model and Fixed-Uncertainty Model) on each

622 contrast-session with 10-fold cross-validation, and reported the aggregated test set scores as was

623 done for the original data.

624 **Data availability** All figures except for Figure 1 were generated from raw data or processed data.

625 The data generated and/or analyzed during the current study are available from the corresponding

626 author upon reasonable request. No publicly available data was used in this study.

627 **Code availability** Code used for modeling and training the deep neural networks as well as for

628 figure generation can be viewed and downloaded from `https://github.com/eywalker/`

629 `v1_likelihood`. All other code used for analysis including data selection and decision model

630 fitting can be found at `https://github.com/eywalker/v1_project`. Finally, code

631 used for elecrophysiology data processing can be found in the Tolias lab GitHub organization

632 `https://github.com/atlab`.

633 **Statistics** Throughout the paper, the level of significance is indicated as * for p <0.05, ** for

634 p <0.01 and *** for p <0.005. Exact p values less than 0.001 were reported as p <0.001. All

635 statistical tests used were two-tailed paired two-sample t-test, unless specified otherwise. Wherever

636 reported, data are means and error bars indicate standard error of the means computed as $\frac{\sigma}{\sqrt{n}}$

637 where $\sigma$ is the standard deviation and $n$ is the size of the sample within the bin, unless specified

638 otherwise.

# References

1. Laplace, P.-S. *Theorie Analytique des Probabilites* (Ve Courcier, Paris, 1812).

2. von Helmholtz, H. Versuch einer erweiterten Anwendung des Fechnerschen Gesetzes im farbensystem. *Z. Psychol. Physiol. Sinnesorg* **2**, 1–30 (1891).

3. Knill, D. C. & Richards, W. (eds.) *Perception As Bayesian Inference* (Cambridge University Press, New York, NY, USA, 1996).

4. Kersten, D., Mamassian, P. & Yuille, A. Object perception as Bayesian inference. *Annual review of psychology* **55**, 271–304 (2004).

5. Knill, D. C. & Pouget, A. The Bayesian brain: the role of uncertainty in neural coding and computation. *Trends in Neurosciences* **27**, 712–719 (2004).

6. Ma, W. J. & Jazayeri, M. Neural Coding of Uncertainty and Probability. *Annual review of neuroscience* **37**, 205–220 (2014).

7. Alais, D. & Burr, D. The ventriloquist effect results from near-optimal bimodal integration. *Current Biology* **14**, 257–262 (2004).

8. Ernst, M. O. & Banks, M. S. Humans integrate visual and haptic information in a statistically optimal fashion. *Nature* **415**, 429–433 (2002). NIHMS150003.

9. Ma, W. J., Beck, J. M., Latham, P. E. & Pouget, A. Bayesian inference with probabilistic population codes. *Nature Neuroscience* **9**, 1432–1438 (2006). NIHMS150003.

10. Beck, J. M. *et al.* Probabilistic Population Codes for Bayesian Decision Making. *Neuron* **60**, 1142–1152 (2008). 1507.01561.

11. Pouget, A., Dayan, P. & Zemel, R. Information processing with population codes. *Nature reviews. Neuroscience* **1**, 125–32 (2000).

12. Pouget, A., Dayan, P. & Zemel, R. S. Inference and Computation with Population Codes. *Annu. Rev. Neurosci* **26**, 381–410 (2003).

13. Ma, W. J., Beck, J. M. & Pouget, A. Spiking networks for Bayesian inference and choice. *Current Opinion in Neurobiology* **18**, 217–222 (2008).

14. Graf, A. B. A., Kohn, A., Jazayeri, M. & Movshon, J. A. Decoding the activity of neuronal populations in macaque primary visual cortex. *Nature Publishing Group* **14**, 239–245 (2011).

15. Yang, T. & Shadlen, M. N. Probabilistic reasoning by neurons. *Nature* **447**, 1075–1080 (2007). NIHMS150003.

16. Fetsch, C. R., Pouget, A., Deangelis, G. C. & Angelaki, D. E. Neural correlates of reliability-based cue weighting during multisensory integration. *Nature Neuroscience* **15**, 146–154 (2012). NIHMS150003.

17. Qamar, A. T. *et al.* Trial-to-trial, uncertainty-based adjustment of decision boundaries in visual categorization. *Proceedings of the National Academy of Sciences* **110**, 20332–20337 (2013).

18. LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436–444 (2015).

19. Goodfellow, Ian, Bengio, Yoshua, Courville, A. Deep Learning. *MIT Press* (2016). `arXiv:` `1312.6184v5`.

20. Seung, H. S. & Sompolinsky, H. Simple models for reading neuronal population codes. *Proc.Natl.Acad.Sci.* **90**, 10749–10753 (1993). `NIHMS150003`.

21. Sanger, T. D. Probability density estimation for the interpretation of neural population codes. *Journal of neurophysiology* **76**, 2790–3 (1996).

22. Zemel, R. S., Dayan, P. & Pouget, A. Probabalistic interpretation of population codes. *Neural Comp.* **10**, 403–430 (1998).

23. Jazayeri, M. & Movshon, J. A. Optimal representation of sensory information by neural populations. *Nature Neuroscience* **9**, 690–696 (2006). `10.1021/nl3012853|NanoLett.` `2012,12,36023608`.

24. Averbeck, B. B. & Lee, D. Effects of Noise Correlations on Information Encoding and Decoding. *J Neurophysiol* **95**, 3633–3644 (2006).

25. Ecker, A. S. *et al.* Decorrelated neuronal firing in coritcal micorcircuits. *Science* **327**, 584–587 (2010).

26. Ecker, A. S., Berens, P., Tolias, A. S. & Bethge, M. The Effect of Noise Correlations in Populations of Diversely Tuned Neurons. *Journal of Neuroscience* **31**, 14272–14283 (2011). `NIHMS150003`.

27. Ecker, A. S. *et al.* State dependence of noise correlations in macaque primary visual cortex. *Neuron* **82**, 235–248 (2014).

28. van Bergen, R. S. & Jehee, J. F. Modeling correlated noise is necessary to decode uncertainty. *NeuroImage* (2017). `1708.04860`.

29. Denfield, G. H., Ecker, A. S., Shinn, T. J., Bethge, M. & Tolias, A. S. Attentional fluctuations induce shared variability in macaque primary visual cortex. *Nature Communications* **9**, 2654 (2018).

30. Ma, W. J. Signal detection theory, uncertainty, and Poisson-like population codes. *Vision Research* **50**, 2308–2319 (2010).

31. Van Bergen, R. S., Ji Ma, W., Pratte, M. S. & Jehee, J. F. Sensory uncertainty decoded from visual cortex predicts behavior. *Nature Neuroscience* **18**, 1728–1730 (2015). `15334406`.

32. Tolhurst, D. J., Movshon, J. A. & Dean, A. F. The statistical reliability of signals in single neurons in cat and monkey visual cortex. *Vision Research* **23**, 775–785 (1983).

33. Shadlen, M. N. & Newsome, W. T. The variable discharge of cortical neurons: implications for connectivity, computation, and information coding. *The Journal of neuroscience : the official journal of the Society for Neuroscience* **18**, 3870–96 (1998).

34. Angelaki, D. E., Humphreys, G. & DeAngelis, G. C. Multisensory Integration. *Journal Of The Theoretical Humanities* **19**, 452–458 (2009).

35. Ma, W. J., Navalpakkam, V., Beck, J. M., van den Berg, R. & Pouget, A. Behavior and neural basis of near-optimal visual search. *Nature Neuroscience* **14**, 783–790 (2011). NIHMS150003.

36. Beck, J. M., Latham, P. E. & Pouget, A. Marginalization in Neural Circuits with Divisive Normalization. *Journal of Neuroscience* **31**, 15310–15319 (2011). NIHMS150003.

37. Ma, W. J. & Rahmati, M. Towards a Neural Implementation of Causal Inference in Cue Combination. *Multisensory Research* **26**, 159–176 (2013).

38. Orhan, A. E. & Ma, W. J. Efficient probabilistic inference in generic neural networks trained with non-probabilistic feedback. *Nature Communications* **8**, 138 (2017).

39. Brainard, D. H. The Psychophysics Toolbox. *Spatial vision* **10**, 433–6 (1997).

40. Tolias, A. S. *et al.* Recording Chronically From the Same Neurons in Awake, Behaving Primates. *Journal of Neurophysiology* **98**, 3780–3790 (2007).

41. Subramaniyan, M., Ecker, A. S., Berens, P. & Tolias, A. S. Macaque Monkeys Perceive the Flash Lag Illusion. *PLoS ONE* **8**, e58788 (2013).

42. Quiroga, R. Q., Nadasdy, Z. & Ben-Shaul, Y. Unsupervised Spike Detection and Sorting with Wavelets and Superparamagnetic Clustering. *Neural Computation* **16**, 1661–1687 (2004).

43. Kohn, A. & Movshon, J. A. Adaptation changes the direction tuning of macaque MT neurons. *Nature Neuroscience* **7**, 764–772 (2004).

44. Richard, M. D. & Lippmann, R. P. Neural Network Classifiers Estimate Bayesian a posteriori Probabilities. *Neural Computation* **3**, 461–483 (1991).

45. Kline, D. M. & Berardi, V. L. Revisiting squared-error and cross-entropy functions for training neural network classifiers. *Neural Computing and Applications* **14**, 310–318 (2005).

46. Kullback, S. & Leibler, R. A. On Information and Sufficiency. *The Annals of Mathematical Statistics* **22**, 79–86 (1951). 1511.00860.

47. MacKay, D. J. C. *Information Theory , Inference , and Learning Algorithms*, vol. 22 (Cambridge University Press, Cambridge, UK, 2003). arXiv:1011.1669v3.

48. Srivastava, N., Hinton, G., Krizhevsky, A. & Salakhutdinov, R. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *Journal of Machine Learning Research* **15**, 1929–1958 (2014).

49. Prechelt, L. Early stopping-but when? In *Neural Networks: Tricks of the Trade, This Book is an Outgrowth of a 1996 NIPS Workshop*, 55–69 (Springer-Verlag, London, UK, UK, 1998).

50. Mamassian, P. & Landy, M. S. Observer biases in the 3D interpretation of line drawings. *Vision Research* **38**, 2817–2832 (1998).

51. Acerbi, L., Vijayakumar, S. & Wolpert, D. M. On the Origins of Suboptimality in Human Probabilistic Inference. *PLoS Computational Biology* **10**, e1003661 (2014).

52. Acerbi, L. & Ma, W. J. Practical Bayesian Optimization for Model Fitting with Bayesian Adaptive Direct Search. In *Advances in Neural Information Processing Systems 30*, 1836–1846 (2017). `1705.04405`.

**Contributions**   All authors designed the experiments and developed the theoretical framework. R.J.C. trained the first monkey, and R.J.C. and E.Y.W. recorded data from this monkey. E.Y.W. trained and recorded from the second monkey. E.Y.W. performed all data analyses. E.Y.W. wrote the manuscript, with contributions from all authors.

**Competing Interests**   The authors declare that they have no competing financial interests.

**Correspondence**   Correspondence and requests for materials should be addressed to A.S.T. (astolias@bcm.edu).