

# A neural basis of probabilistic computation in visual cortex

Edgar Y. Walker,<sup>1,2†</sup> R. James Cotton,<sup>1,2,3†</sup> Wei Ji Ma,<sup>4‡</sup>  
Andreas S. Tolias<sup>1,2,5‡\*</sup>

<sup>1</sup>Center for Neuroscience and Artificial Intelligence,  
Baylor College of Medicine, TX, USA

<sup>2</sup>Department of Neuroscience, Baylor College of Medicine, TX, USA

<sup>3</sup>Now at: Shirley Ryan AbilityLab, IL, USA

<sup>4</sup>Center for Neural Science and Department of Psychology,  
New York University, NY, USA

<sup>5</sup>Department of Electrical and Computer Engineering,  
Rice University, TX, USA

\*To whom correspondence should be addressed; E-mail: astolias@bcm.edu.

†‡These authors contributed equally to this work

## Abstract:

Bayesian models of behavior suggest that organisms represent uncertainty associated with sensory variables. However, 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. We studied the neural code of uncertainty by simultaneously recording population activity from the primate visual cortex during a visual categorization task in which trial-to-trial uncertainty about stimulus orientation was relevant for

**the decision. We decoded the likelihood function from the trial-to-trial population activity and found that it predicted decisions better than a point estimate of orientation. This remained true when we conditioned on the true orientation, suggesting that internal fluctuations in neural activity drive behaviorally meaningful variations in the likelihood function. Our results establish the role of population-encoded likelihood functions in mediating behavior, and provide a neural underpinning 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<sup>1</sup> and von Helmholtz<sup>2</sup>—has been widely used to explain human and animal perception<sup>3–6</sup>. At its core 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 world state variables (e.g. direction of motion of a stimulus). In spite of the prevalence of Bayesian theories in neuroscience, evidence to support them stems primarily from behavioral studies (e.g.<sup>7,8</sup>). Consequently, the manner in which probability distributions are encoded in the brain remains unclear, and, thus, the neural code of uncertainty is unknown.

It has been hypothesized that a critical feature of the neural code of uncertainty, which is shared throughout the sensory processing chain in the neocortex, is that the same neurons that encode a specific world state variable (e.g. stimulus orientation in V1) also encode the uncertainty about that variable (Fig. 1a). Therefore neurons multiplex both a point estimate of a sensory variable and the associated uncertainty about it<sup>9,10</sup>. Specifically, according to the probabilistic population coding (PPC) hypothesis<sup>9,10</sup>, inference in the brain is performed by inverting a generative model of neural population activity. Under this coding scheme, neural populations in V1, for example, that encode stimulus orientation also encode the associated

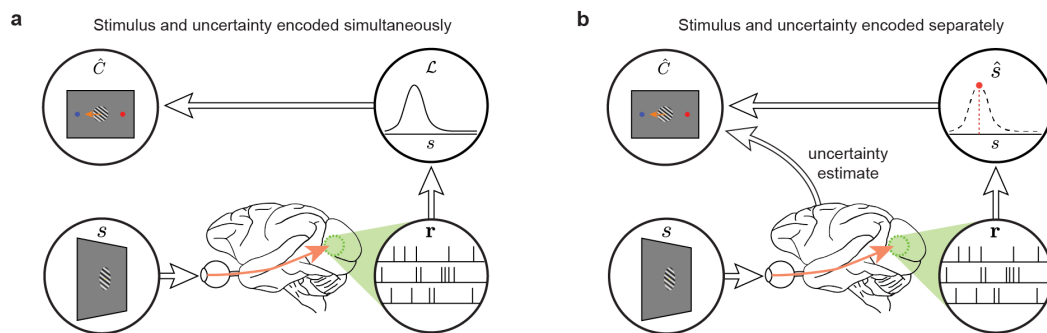


Figure 1: Alternative models of uncertainty information encoding. **a**, The recorded cortical population  $r$  responding to sensory stimulus  $s$  encodes stimulus estimate and uncertainty simultaneously in the form of likelihood function  $\mathcal{L}$  which is subsequently used in making a decision  $\hat{C}$  as the subject performs a visual classification task. **b**, The recorded cortical population only encodes a point estimate of the stimulus  $\hat{s}$  while an estimate of the sensory uncertainty is made by other (unrecorded) cortical populations. The information is subsequently combined to lead to the subject's decision  $\hat{C}$ .

uncertainty in the form of the sensory likelihood function—the probability of observing a given pattern of neural activity across hypothesized stimulus values<sup>9,11,12</sup>. The form of the likelihood function is related to the probability distribution describing neural variability (“noise”) for a given stimulus. A sensory likelihood function is often unimodal<sup>13,14</sup>, and its width could in principle serve as a measure of the sensory uncertainty about the stimulus. Whether the brain uses this particular uncertainty quantity in its decisions is unknown. Alternatively, it may be the case that the neural population that encodes an estimate of a sensory variable (e.g. stimulus orientation in V1) does not carry information about the associated uncertainty (Fig. 1b).

We recorded the activity of V1 cortical populations while monkeys performed a visual classification task in which the trial-by-trial uncertainty information is beneficial to the animal<sup>15</sup>. To decode the trial-by-trial likelihood functions from the V1 population responses, we developed a novel technique based on deep learning<sup>16,17</sup>. Importantly, we performed all analyses conditioned on the contrast—an overt driver of uncertainty—and performed further orientation-

conditioned analyses to isolate the effect of random fluctuations in the decoded likelihood function on behavior. We found that using the trial-to-trial changes in the shape of the likelihood function allowed us to better predict the behavior than using a likelihood function with a fixed shape shifted by a point estimate. Therefore, we provide the first evidence that in perceptual decision-making, the same cortical population that encodes a sensory variable also encodes its trial-by-trial sensory uncertainty information, which is used to mediate perceptual decisions, consistent with the theory of PPC.

## Results

### Behavior

Two Rhesus macaques (*Macacca mulatta*) were trained on an orientation classification task designed such that the optimal performance required the use of trial-by-trial uncertainty. On each trial, one of two stimulus classes ( $C = 1$  or  $C = 2$ ) was chosen at random with equal probability. Each class was defined by a Gaussian probability distribution over the orientation. The two distributions shared the same mean but had different standard deviations (Fig. 2a). An orientation was drawn from the distribution belonging to the selected class, and a drifting grating stimulus with that orientation was then presented to the animal (Fig. 2b). In a given recording session, at least three distinct contrasts were selected at the beginning of the session, and on each trial, one of these values was randomly selected.

In our previous study<sup>15</sup>, we designed this task so that an optimal Bayesian observer would incorporate the trial-by-trial sensory uncertainty about stimulus orientation in making classification decisions. Indeed, decisions of both humans and monkeys seemed to utilize trial-by-trial uncertainty about the stimulus orientation. In the current study, one of the two monkeys (Monkey L) was the same monkey that participated in the previous study and thus has been shown to have learned the task well. A second monkey (Monkey T) was also trained on the task and

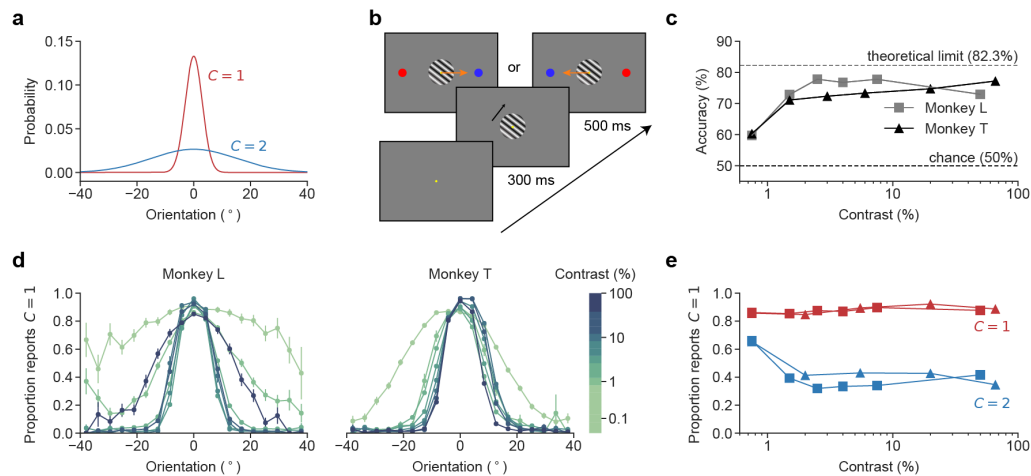


Figure 2: Behavioral task. **a**, The stimulus orientation distributions for the two classes. The two distributions shared the same mean ( $\mu = 0^\circ$ ) but differed in their standard deviations ( $\sigma_1 = 3^\circ$  and  $\sigma_2 = 15^\circ$ ). **b**, Time course of a single trial. The subject 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 was 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**.

74 closely matched the performance of Monkey L (Fig. 2c). Both animals had psychometric curves  
75 displaying the expected strong dependence on both contrast and orientation (Fig. 2d,e).

76 In our analyses, we grouped the trials with the same contrast within the same session and  
77 refer to such a group as a “contrast-session”.

## 78 **Decoding likelihood function from V1**

79 Each monkey was implanted with a chronic multi-electrode (Utah) array in the parafoveal pri-  
80 mary visual cortex (V1) to record the simultaneous cortical population activity as the subjects  
81 performed the orientation classification task (Fig. 3a).

82 A total of 61 and 71 sessions were analyzed from Monkey L and Monkey T for a total of  
83 110,695 and 192,631 trials, respectively (Supplementary Fig. 1). In each recording session,  
84 up to 96 channels were recorded. On each trial and for each channel, we computed the to-  
85 tal number of spikes that occurred during the 500 ms of stimulus presentation preceding the  
86 decision-making cue (Fig. 3a), yielding a vector of population responses  $\mathbf{r}$  used in the subse-  
87 quent analyses (Fig. 3b).

88 Existing computational methods for decoding the trial-by-trial likelihood function from the  
89 cortical population activities typically make strong parametric assumptions about the stimulus  
90 conditioned distribution of the population response (i.e. the generative model of the population  
91 response). For example, population responses to a stimulus can be modeled as an independent  
92 Poisson distribution, allowing each recorded unit to be characterized by a simple tuning curve  
93 (which may be further parameterized)<sup>14,18–22</sup>. While this simplifying assumption makes com-  
94 puting the trial-by-trial likelihood function straightforward, disregarding potential correlations  
95 among the units in population responses (i.e. noise correlations and internal brain state fluctu-  
96 ations<sup>23–28</sup>) can lead to biased estimates of the likelihood function and limits the generality of  
97 this approach. While more generic parametric models—such as Poisson-like distributions—of

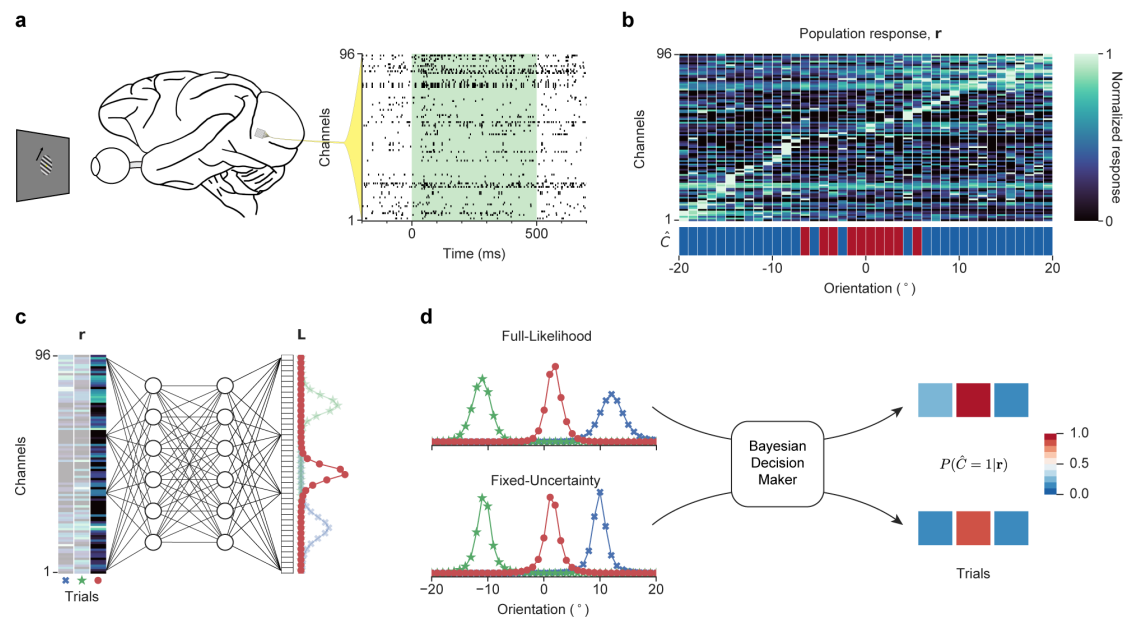


Figure 3: Encoding and decoding of the stimulus orientation. **a**, An example of 96 channels spike traces from a single trial (Monkey T). The vector of spike counts,  $r$ , was accumulated over the pre-saccade stimulus presentation period (time 0-500 ms, green shade). **b**, The population response for the selected trials from a single contrast-session (Monkey T, 64% contrast). Column: a population response  $r$  on a trial randomly drawn from the trials falling into a specific orientation bin. Row: 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. Subject's class decision is indicated by red and blue color patches for  $\hat{C} = 1$  and  $\hat{C} = 2$ , respectively. **c**, A schematic of a DNN for the Full-Likelihood decoder, mapping  $r$  to the decoded likelihood function  $L$ . All likelihood functions are area-normalized. **d**, Two models of likelihood decoder  $M$ . In the *Full-Likelihood decoder*, the likelihood  $L$  was decoded without any constraints on the shape. In the *Fixed-Uncertainty decoder*, all decoded likelihood functions shared the same shape but differed in the location of the center based on the population response.

For both decoders, the resulting likelihood functions were fed into parameterized Bayesian decision models to yield the decision prediction  $p(\hat{C} = 1|r, M)$ .

population distributions have been proposed<sup>9,10,15,29,30</sup>, they still impose restrictive assumptions.

We devised a technique based on deep learning to decode the trial-by-trial likelihood function from the V1 population response. This neural network-based likelihood decoder allows us to approximate the information that can be extracted about the stimulus orientation from the cortical population responses. The network was *not* used as a model of how the rest of the brain extracts and processes the information present in the population, but rather to decode it and demonstrate that it is used behaviorally.

We trained a fully connected deep neural network (DNN)<sup>17</sup> to predict the per-trial likelihood function  $\mathcal{L}(\theta) \equiv p(\mathbf{r}|\theta)$  over stimulus orientation  $\theta$  from the vectorized population response  $\mathbf{r}$  (Fig. 3c; for details on the network architecture, training objective, and hyperparameter selection see Methods and Supplementary Table 1). A separate network was trained for each contrast-session and no behavioral data were utilized in training the DNN.

Using a DNN to decode the likelihood function avoids the restrictive parametric assumptions described above and provides a strictly more flexible method, often capturing decoding under known distributions as a special case (Supplementary Fig. 2). We demonstrate this by showing the DNN can recover the ground-truth likelihood function from simulated responses sampled from known distributions (Supplementary Fig. 3; refer to Methods for the simulation details).

The likelihood functions decoded by the DNNs exhibited the expected dependencies on the overt drivers of uncertainty such as contrast (Fig. 4a-c): the width of the likelihood function is higher at lower contrast (Fig. 4d).

## **Trial to trial uncertainty improves behavioral predictions**

To assess whether the uncertainty decoded from population responses in the form of sensory likelihood functions mediate the behavioral outcome (perceptual decision) as we hypothesized, it is critical that we appropriately condition the analysis on the stimulus. To illustrate the impor-



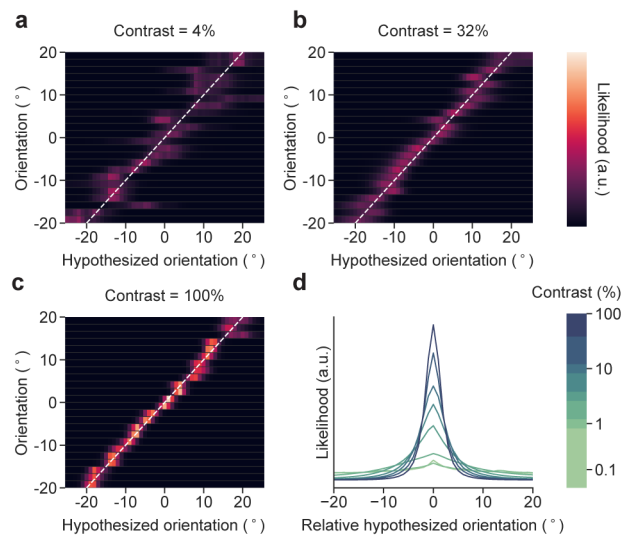


Figure 4: Likelihood functions decoded by the trained neural networks. **a-c**, Example decoded likelihood functions from three contrast-sessions from Monkey T. Each row represents the decoded likelihood function over the hypothesized orientation for a randomly selected trial within the specific orientation 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 occurred at  $0^\circ$ . The centered likelihood functions were then averaged across all trials within the same contrast bin.

tance of conditioning on the stimulus to determine if the decoded likelihood function mediates perceptual decisions, consider a typical perceptual decision-making task (like ours) (Supplementary Fig. 4) where the subject views a stimulus  $s$ , which elicits a population response  $\mathbf{r}$ , for example in V1. Here, by “stimulus”, we refer collectively to all aspects of a visual stimulus, such as its contrast and orientation. Stimulus information is eventually relayed to decision-making areas (e.g. prefrontal cortex), leading the animal to make a classification decision  $\hat{C}$ . We decode the likelihood function  $\mathcal{L}$  from the recorded population activity  $\mathbf{r}$ . Because variation in the stimulus (e.g. orientation or contrast) across trials can drive variation both in the decoded likelihood function and in the animal’s decision, one may find a significant relationship between  $\mathcal{L}$  and  $\hat{C}$ , even if the likelihood function estimated from the recorded population  $\mathbf{r}$  does not mediate the decision. When the stimulus is fixed, random fluctuations in the population response  $\mathbf{r}$  can still result in variations in  $\mathcal{L}$ . If the likelihood function truly mediates the decision, we expect that such variation in  $\mathcal{L}$  would account for variation in  $\hat{C}$ . Therefore, to demonstrate that the likelihood  $\mathcal{L}$  mediates the decision  $\hat{C}$ , it is imperative to show a correlation between  $\mathcal{L}$  and  $\hat{C}$  conditioned on the stimulus  $s$ .

As we varied the stimulus contrast from trial to trial in our task, the expected uncertainty about the stimulus orientation varied, and one would expect the monkeys to represent and make use of their trial-by-trial sensory uncertainty in making decisions. However, we make a much stronger claim here: even at a fixed contrast, because of random fluctuations in the population response<sup>31,32</sup>, we predict (1) the uncertainty encoded in the population, that is, the likelihood function, will still fluctuate from trial to trial, and (2) the effect of such fluctuations will manifest in the monkey’s decisions on a trial-by-trial basis.

We tested this prediction by fitting, separately for each contrast-session, the following two decision models and comparing their performance in predicting the monkey’s trial-by-trial decisions: (1) a Full-Likelihood Model, which utilizes the trial-by-trial uncertainty information

decoded from the population response in the form of the likelihood function obtained from the neural-network based likelihood decoder (Full-Likelihood decoder) described above (Fig. 3d), and (2) a Fixed-Uncertainty Model, which utilizes an alternative neural-network based likelihood decoder (Fixed-Uncertainty decoder) that learns a single, fixed-shape likelihood function whose location is shifted from trial to trial based on the population response (Supplementary Fig. 5). The Fixed-Uncertainty Model captures the alternative hypothesis in which the recorded sensory population only encodes a point estimate of the sensory variable (i.e. mean of the likelihood function) and the estimate of the sensory uncertainty is encoded elsewhere, signified by the fixed shape of the likelihood function fitted for each contrast-session under this model (Fig. 1b). Generally, the likelihood function decoded by Fixed-Uncertainty decoder closely approximated the likelihood function decoded by the Full-Likelihood decoder (Supplementary Fig. 5). We use the term decoder for the DNN that returns estimated likelihood functions, and the term decision model for the mapping from likelihood function to decision.)

In both models, the decoded likelihood functions were fed into the Bayesian decision maker to yield trial-by-trial predictions of the subject's decision in the form of  $p(\hat{C}|\mathbf{r}, M)$ , or the likelihood of subject's decisions  $\hat{C}$  conditioned on the population response  $\mathbf{r}$  and the decision model  $M$ . The Bayesian decision maker computed the posterior probability of each class and used these to produce a stochastic decision. The means of the class distributions assumed by the observer, the class priors, the lapse rate, and a parameter to adjust the exact decision-making strategy were used as free parameters (Supplementary Fig. 6, refer to Methods for details). The model parameters were fitted by maximizing the total log likelihood over all trials for each contrast-session  $\sum_i \log p(\hat{C}_i|\mathbf{r}_i, M)$ . The fitness of the models was assessed through cross-validation, and we reported mean and total log likelihood of the models across all trials in the test set.

Both models incorporated trial-by-trial changes in the point estimate of the stimulus orien-

tation (e.g. the mean of the likelihood function) and only differed in whether they contained additional uncertainty information about the stimulus orientation carried by the trial-by-trial fluctuations in the shape of the likelihood function decoded from the same population that encoded the point estimate. We use the term “shape” to refer to all aspects of the likelihood function besides its mean, including its width. If the fluctuations in the shape of the likelihood function truly captured the fluctuations in the sensory uncertainty as represented and utilized by the animal, one would expect the Full-Likelihood Model to yield better trial-by-trial predictions of the monkey’s decisions than the Fixed-Uncertainty Model.

We observed that both models predicted the monkey’s behavior well across all contrasts (Supplementary Fig. 7), reaching up to 90% accuracy. We also observed that the performance of the decision models using likelihood functions that were decoded by the neural networks was superior to the models using likelihood functions that were decoded with more traditional parametric generative models (independent Poisson distribution and Poisson-like distribution) (Supplementary Fig. 8; refer to Methods for details). The Full-Likelihood Model consistently outperformed the Fixed-Uncertainty Model across contrasts and for both monkeys (Fig. 5a,b; trial log likelihood differences between the Full-Likelihood and Fixed-Uncertainty Model: Monkey L: paired t-test,  $t(110694) = 11.06$ ,  $p < 0.001$ ,  $\delta_{\text{total}} = 11.0 \times 10^2$  and Monkey T:  $t(192610) = 11.03$ ,  $p < 0.001$ ,  $\delta_{\text{total}} = 11.3 \times 10^2$ ;  $\delta_{\text{total}}$  is the total log likelihood difference across all trials). This result shows that the trial-by-trial fluctuations in the shape of the likelihood function are informative about the monkey’s trial-by-trial decisions, demonstrating that decision-relevant sensory uncertainty information is contained in population responses that can be captured by the shape of the full likelihood function. This finding in turn strongly supports the hypothesis that visual cortex encodes stimulus uncertainty through the shape of the full likelihood function on a trial-by-trial basis.

We repeated this analysis after splitting the data into the first and second 250ms of stimulus

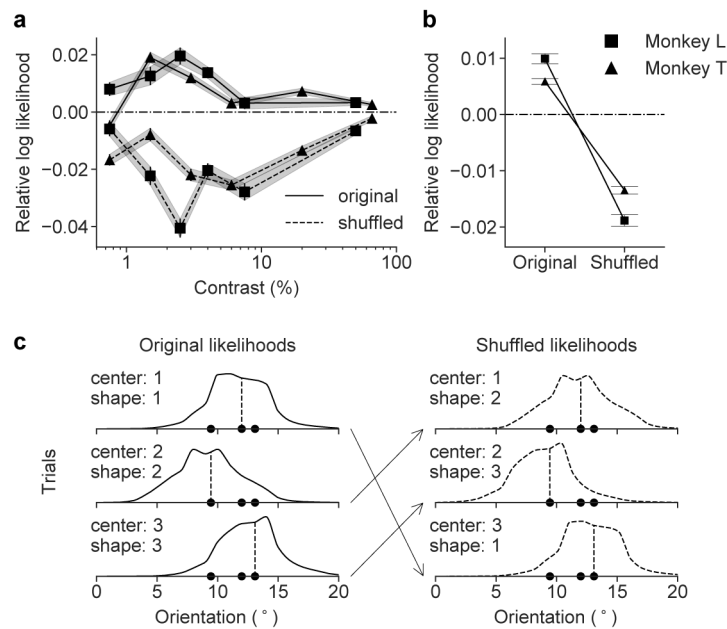


Figure 5: Model performance. **a**, Average trial-by-trial performance of the Full-Likelihood Model relative to the Fixed-Uncertainty Model across contrasts, measured as the average trial difference in the log likelihood. The results for the original (unshuffled) and the shuffled data are shown in solid and dashed lines, respectively. The squares and triangles mark Monkey L and T, respectively. **b**, Relative model performance summarized across all contrasts. Performance on the original and the shuffled data is shown individually for both monkeys. The difference between the Full-Likelihood and Fixed-Uncertainty Models was significant with  $p < 0.001$  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 ( $p < 0.001$  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.

presentation. We found a similar improvement for the Full-Likelihood model over the Fixed-Uncertainty model in both periods (Supplementary Fig. 9).

We next asked how meaningful our effect sizes (model performance differences) are. To answer this question, we simulated the monkey's responses across all trials and contrast-sessions taking the trained Full-Likelihood Model to be the ground truth, and then retrained the Bayesian decision makers in the Full-Likelihood Model and the Fixed-Uncertainty Model from scratch on the simulated data. This approach yields a theoretical upper bound on the observable difference between the two models if the Full-Likelihood Model was the true model of the monkeys' decision-making process.

We observed that the expected total upper bound log likelihood differences between the Full-Likelihood Model and the Fixed-Uncertainty Model of  $(37.1 \pm 1.5) \times 10^2$  and  $(36.0 \pm 1.3) \times 10^2$  based on the simulations (representing mean  $\pm$  standard deviation across 5 repetitions of simulation for Monkey L and T, respectively) were larger but in the same order of the magnitude as the observed model performance differences ( $11.0 \times 10^2$  and  $11.3 \times 10^2$  total log likelihood differences across all trials for Monkey L and T, respectively), suggesting that our effect sizes are meaningful and that the Full-Likelihood Model is a reasonable approximate description of the monkey's true decision-making process (Supplementary Fig. 10).

## Stimulus dependent changes in uncertainty

We observed that for some contrast-sessions, the average width of the likelihood function showed a dependence on the stimulus orientation (Supplementary Figure. 11). By design, the Fixed-Uncertainty Model cannot capture this stimulus dependent change in uncertainty, which could contribute to it under-performing the Full-Likelihood Model (Supplementary Fig. 4).

To rule this out, we shuffled the shapes of the decoded likelihood functions across trials within the same orientation bin, separately for each contrast-session. This shuffling preserves

the average stimulus-dependent change in uncertainty and trial-by-trial correlation between the mean of the likelihood function and the decision (Fig. 5c), while removing the trial-by-trial correlation between the shape of the likelihood function and the behavioral decision conditioned on the stimulus orientation.

By design, the Fixed-Uncertainty Model makes identical predictions on the original and the shuffled data. If the Full-Likelihood Model outperformed the Fixed-Uncertainty Model simply because it captured spurious correlations between the stimulus orientation and the shape of the likelihood function, then it should outperform the Fixed-Uncertainty model by the same amount on the shuffled data. However, if the better behavioral predictions come from the trial-by-trial fluctuations in the likelihood shape as we hypothesized, one would expect this difference to disappear on the shuffled data. Indeed, the shuffling of the likelihood function shapes abolished the improvement in prediction performance that the Full-Likelihood Model had over the Fixed-Uncertainty Model. In fact, the Full-Likelihood Model consistently underperformed the Fixed-Uncertainty Model on the shuffled data (Fig. 5a,b; trial log likelihood difference between the Full-Likelihood Model and the Fixed-Uncertainty Model on the shuffled data: Monkey L: paired t-test  $t(110694) = -18.44$ ,  $p < 0.001$ ,  $\delta_{\text{total}} = -20.9 \times 10^2$  and Monkey T:  $t(192610) = -20.15$ ,  $p < 0.001$ ,  $\delta_{\text{total}} = -25.9 \times 10^2$ ;  $\delta_{\text{total}}$  is the total log likelihood difference across all trials). Therefore, there were significant performance differences in Full-Likelihood Model between the unshuffled and shuffled data (trial log likelihood difference: Monkey L: paired t-test  $t(110694) = 33.34$ ,  $p < 0.001$ ,  $\delta_{\text{total}} = 31.9 \times 10^2$  and Monkey T:  $t(192610) = 34.52$ ,  $p < 0.001$ ,  $\delta_{\text{total}} = 37.2 \times 10^2$ ).

To confirm our effect sizes were appropriate, we again compared these values to those obtained from simulations in which we took the Full-Likelihood Model to be the ground truth (Supplementary Fig. 10). The simulations yielded total log likelihood differences of the Full-Likelihood Model between the unshuffled and shuffled data of  $(36.2 \pm 2.2) \times 10^2$  (Monkey L)

and  $(40.7 \pm 1.5) \times 10^2$  (Monkey T) (mean  $\pm$  standard deviation across 5 repetitions), similar in magnitude to the observed values.

Taken together, the shuffling analyses show that the better prediction performance of the Full-Likelihood Model is not due to the confound between the stimulus and the likelihood shape. We conclude that the trial-by-trial likelihood function decoded from the population represents behaviorally relevant stimulus uncertainty information, even when conditioned on the stimulus.

## Attribution analysis

To assess whether the same population encoding the best point estimate (i.e. mean of the likelihood function) also encoded the uncertainty regarding that estimate (i.e. shape of the likelihood function), as we hypothesized to be the case, we performed attribution analysis<sup>33</sup> on the trained Full-Likelihood decoder. Through this analysis, we ask how much of the changes in either (1) the mean of the likelihood  $\mu_L$  (i.e. surrogate for the best point estimate) or (2) the standard deviation of the likelihood function  $\sigma_L$  (i.e. surrogate measure of the uncertainty) can be attributed back to each input multiunit, yielding attribution  $A_\mu$  and  $A_\sigma$ , respectively. The question of feature attribution is a very active field of research in machine learning, and multiple methods of attribution computation exist<sup>33–35</sup>. Here we have selected three different methods of computing attribution scores: saliency maps<sup>34</sup>, gradient  $\times$  input<sup>33</sup>, and DeepLIFT<sup>35</sup> (refer to Methods for the details of attribution computation).

We observed that across all three attribution methods, multiunits with high  $\mu_L$  attribution tended to have high  $\sigma_L$  attribution, and vice versa, giving rise to high degree of correlation between  $A_\mu$  and  $A_\sigma$  (Fig. 6a). If distinct subpopulations were involved in encoding the point estimate and the uncertainty as found in the likelihood function, we would have expected multiunits with a high  $\mu_L$  attribution score to have a low  $\sigma_L$  attribution score, and vice versa, there-



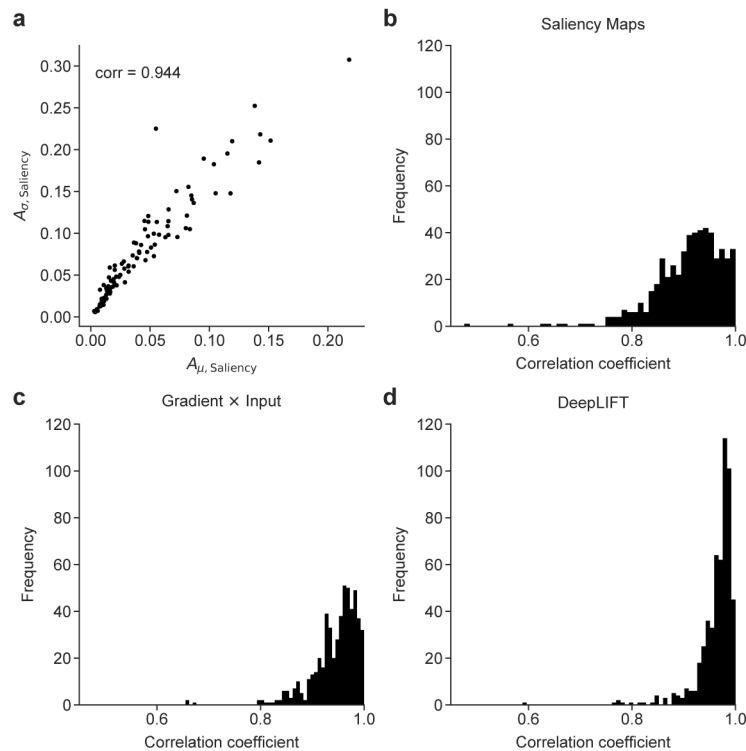


Figure 6: Attribution analysis for means and standard deviations of the likelihood functions. **a**, Attribution of 96 input multiunits to the likelihood mean  $A_{\mu, \text{Saliency}}$  vs. standard deviation  $A_{\sigma, \text{Saliency}}$  computed based on saliency maps for an example contrast session (Monkey T, 32% contrast). **b-d**, Distribution of correlation coefficients between  $A_{\mu}$  and  $A_{\sigma}$  for multi units across all contrast sessions for both monkeys, computed based on different attribution methods.

fore leading to negative correlation between  $A_{\mu}$  and  $A_{\sigma}$ . However, we observed that across all contrast-sessions from both monkeys,  $A_{\mu}$  was strongly positively correlated with  $A_{\sigma}$  regardless of the exact attribution method used, suggesting that the highly overlapping subpopulations are involved in encoding both the point estimate and the uncertainty of the likelihood function, as we hypothesized would be the case (Fig. 6b-d).

## Discussion

Given the stochastic nature of the brain, repeated presentations of identical stimuli elicit variable responses. The covariation between neuronal activity fluctuations and perceptual choice has been studied extensively at the level of single neurons, originating with the pioneering work of Campbell & Kulikowski<sup>36</sup> and Britten et al.<sup>37</sup>. Here, we go beyond this literature by examining the hypothesis that the brain takes into account knowledge of the form of neural variability in order to build a belief over the stimulus of interest on each trial. This belief is captured by the likelihood function and the associated sensory uncertainty, both of which vary from trial to trial with the neural activity. To test this hypothesis, we decoded trial-to-trial likelihood functions from the population activity in visual cortex and used them in conjunction with a highly constrained, theoretically motivated decision model (the Bayesian model) to predict behavior. We found that a model utilizing the full likelihood function predicted the monkeys' choices better than alternative models that ignore variations in the shape of the likelihood function. Our results provide the first population-level evidence in support of the theoretical framework of probabilistic population coding, where the same neurons that encode specific world state variables also encode the uncertainty about those variables. Importantly, under this framework the brain performs Bayesian inference under a generative model of the neural activity.

Our findings were made possible by recording from a large population simultaneously and by using a task in which uncertainty is relevant to the animal. In addition, we decoded likelihood functions using a deep neural network that does not rely on the strong parametric assumptions about the underlying generative model of the population that have dominated previous work. Importantly, we conditioned our analyses on the stimulus to rule out a confounding effect of the stimulus on the observed relationship between the decoded likelihood function and the subject's decision. This approach is critical because previous behavioral studies on cue combination and

Bayesian integration, for instance, always relied on varying stimulus features (e.g. contrast, blur, motion coherence) to manipulate uncertainty<sup>7,8,22,38</sup>. As a result, these studies cannot rule out that any observed correlation between a proposed method of encoding uncertainty and a subject's behavior may be confounded by the stimulus (Supplementary Fig. 4), and they therefore fail to provide a sufficiently rigorous assessment on the representation of uncertainty. Carefully controlling for the effect of stimulus fluctuations allowed us to present rigorous evidence that the trial-by-trial fluctuations in the likelihood functions carry behaviorally relevant stimulus uncertainty information.

After showing that this likelihood function is used behaviorally, what more can we say about the neural encoding of perceptual uncertainty? First, our network learns the *log-likelihood of  $s$* , i.e.  $\log \mathcal{L}(s) = \log p(\mathbf{r}|s) + b(\mathbf{r})$  as a function of  $s$ . We never commit to a particular generative model  $p(\mathbf{r}|s)$  as a function of  $\mathbf{r}$ , as the DNN has an arbitrary offset as a function of  $\mathbf{r}$  (Eq. 1 in Methods). Second, we had to move away from Poisson-like variability to better characterize the responses at the cost of analytic forms and easy interpretability. We see this as a necessary evil; namely, we have shown that making the Poisson-like assumption leads to worse predictions of behavior. That being said, the DNN extends what we know about generative models in visual cortex (e.g. tuning curves, contrast gain); in particular, it allows for rich correlation among units in the population. Third, we would like to stress that we do not believe that the DNN that we use to decode the likelihood is literally implemented in the brain. It remains an important question, and avenue for future research, what kind of transformation, if any, the brain performs in order to utilize and compute with this information.

While the sensory likelihood function is a crucial building block for probabilistic computation in the brain, fundamental questions remain regarding the nature of such computation. First, how do downstream areas process the information contained in sensory likelihood functions to make better decisions? Previous work has manually constructed neural networks for down-

stream computation that relied heavily on the assumption of Poisson-like variability<sup>9,10,15,39–41</sup>. However, more recent work has demonstrated that training generic shallow networks accomplishes the same goal without the need for task-specific manual construction<sup>42</sup>. Second, does each area in a feedforward chain of computation encode a likelihood function over its own variable, with the computation propagating the uncertainty information from one variable to the next? For example, in our task, it is conceivable that prefrontal cortex encodes a likelihood function over class that is derived from a likelihood function over orientation coming in from V1. Third, what are the relative contributions of feedforward, recurrent, and feedback connections to the trial-to-trial population activity and the resulting decoded likelihood functions? Some work has argued strongly for a role of feedback<sup>28,43,44</sup>; in the present work, we are agnostic to this issue. While answering these questions will require major efforts, we expect that our findings will help put those efforts on a more solid footing. In the meantime, our results elevate the standing of Bayesian models of perception from frameworks to describe optimal input-response mappings<sup>45,46</sup> to process models whose internal building blocks—likelihood functions and probability distributions—are more concretely instantiated in neuronal activity<sup>6,47,48</sup>.

## 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: 3.89 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<sup>49</sup>. 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  $0^\circ$  (grating drifting horizontally rightward, positive orientation corresponding to counter-clockwise rotation), but their standard deviations differed:  $\sigma_1 = 3^\circ$  for class 1 ( $C = 1$ ) and  $\sigma_2 = 15^\circ$  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<sup>15</sup>.

Each trial proceeded as follows. A trial was initiated by a beeping sound and the appearance of a fixation target ( $0.15^\circ$  visual angle) in the center of the screen. The monkey fixated on the fixation target for 300 ms within  $0.5^\circ$ – $1^\circ$  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  $4.29^\circ$  from the center with the target diameter of  $0.71^\circ$  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. For illustrative clarity, we used blue to indicate  $C = 2$  throughout this document. The monkey received a juice reward for each correct response (0.10–0.15 mL).

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

task, where the grating stimulus was colored to match the correct class  $C$  for that trial (red for  $C = 1$  and green for  $C = 2$ ). Under this arrangement, the monkey simply learned to saccade to the target matching the color of the grating stimulus, although the grating stimulus orientation information was always present. As the training proceeded, we gradually removed the color from the stimulus, encouraging the monkey to make use of the orientation information in the stimulus to perform the task. Eventually, the color was completely removed, and at that point the monkey was performing the full version of the task.

## Surgical Methods

Our surgical procedures followed a previously established approach<sup>28,50,51</sup>. Briefly, a custom-built titanium cranial headpost was first implanted for head stabilization under general anesthesia using aseptic conditions in a dedicated operating room. After premedication with Dexamethasone (0.25-0.5 mg/kg; 48 h, 24 h and on the day of the procedure) and atropine (0.05 mg/kg prior to sedation), animals were sedated with a mixture of ketamine (10 mg/kg) and xylazine (0.5 mg/kg). During the surgery, anesthesia was maintained using isoflurane (0.5-2%). After the monkey was fully trained, we implanted a 96-electrode microelectrode array (Utah array, Blackrock Microsystems, Salt Lake City, UT, USA) with a shaft length of 1 mm over parafoveal area V1 on the right hemisphere. This surgery was performed under identical conditions as described for headpost implantation. To ameliorate pain, analgesics were given for 7 days following 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<sup>26,28,50</sup>. 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<sup>52</sup>, 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 electrophysiology recordings contaminated by noise in the recording devices (e.g. poor grounding connector resulting in movement noise) or equipment failures. To do so, we established the following trial inclusion criteria:

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



range IQR as follows:  $\sigma_{\text{adj}} = \frac{\text{IQR}}{1.35}$ .

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

We only included trials that satisfied both of the above criteria in our analysis. Empirically, we found the above criteria to be effective in catching obvious anomalies in the spike data while introducing minimal bias into the data. After the application of the criteria, we were left with 110,695 and 192,631 trials for Monkey L and T, thus retaining 98.77% and 99.48% of the total trials, respectively. While this selection criteria allowed us to remove apparent anomaly in the data, we found that the main findings described in this paper were not sensitive to the precise definition of the inclusion criteria.

During each recording session, stimuli were presented under three or more contrast values. In all analyses to follow, we studied the trials from distinct contrast separately for each recording session, 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  $0.29^\circ$  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.

## Full-Likelihood decoder

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<sup>19</sup>. Under this approach, the expected firing rates (i.e. tuning curves) of the  $i^{\text{th}}$  neuron  $E[r_i|\theta] = f_i(\theta)$  must be approximated as well, for example by fitting a parametric function (e.g. von Mises tuning curves<sup>53</sup>) or employing kernel regression<sup>19</sup>. 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  $\mathbf{r}$ .

To this end, we employed a deep neural network (DNN)<sup>16</sup> to directly approximate the likelihood function over the stimulus orientation,  $\theta$ , from the recorded population response  $\mathbf{r}$ . Here we present a brief derivation that serves as the basis of the network design and training objective. Let us assume that  $m$  multiunits were recorded simultaneously in a single recording session, so that  $\mathbf{r} \in \mathbb{R}^m$ . To make the problem tractable, we bin the stimulus orientation  $\theta$  into  $n$  distinct values,  $\theta_1$  to  $\theta_n$  (the derivation holds in general for arbitrarily fine binning of the orientation). With this, the likelihood function can be captured by a vector  $\mathbf{L} \in \mathbb{R}^n$  where  $L_i = \mathcal{L}(\theta_i)$ . Let us assume that we can train some DNN to learn a mapping  $f$  from the population response  $\mathbf{r}$  to the log of the likelihood function  $\mathbf{L}$  up to a constant offset  $b$ . That is,  $f: \mathbb{R}^m \mapsto \mathbb{R}^n$ ,

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

479

480 for some scalar function  $b \in \mathbb{R}$ . As the experimenter, we know the distribution of the stimulus  
481 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  
482 over stimulus orientation  $\theta$  up to some scalar value  $b'(\mathbf{r})$ ,

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

483

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

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

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

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

486

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

488 The goal then is to train the DNN  $f(\mathbf{r})$  such that the overall function  $\mathbf{q}(\mathbf{r})$  matches the  
489 posterior 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  
490 in turn allows the network output  $f(\mathbf{r})$  to approach the log of the likelihood function  $\mathbf{L}$ , up to  
491 a constant  $b(\mathbf{r})$ . For 1-out-of- $n$  classification problems, minimizing the cross-entropy between  
492  $\mathbf{q}(\mathbf{r})$  and the stimulus orientation  $\theta$  for a given  $\mathbf{r}$  lets the overall function  $\mathbf{q}(\mathbf{r})$  approach the true  
493 posterior  $\mathbf{p}(\mathbf{r})$ , as desired<sup>54,55</sup>. To show this, let us start by minimizing the difference between  
494 the model estimated posterior  $\mathbf{q}(\mathbf{r})$  and the true posterior  $\mathbf{p}(\mathbf{r})$  over the distribution of  $\mathbf{r}$ . We do  
495 this by minimizing the loss  $L$  defined as the expected value of the Kullback-Leibler divergence<sup>56</sup>

496 between the two posteriors:

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

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

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

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

497

498 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  
 499 parameters in the network, and  $H(\theta|\mathbf{r})$  is the conditional entropy of  $\theta$  conditioned on  $\mathbf{r}$ , which  
 500 is an unknown but a fixed quantity with respect to  $W$  and the data distribution. Here we used  
 501 the notation  $\mathbf{q}(\mathbf{r}, W)$  to highlight the dependence of the network estimated posterior  $\mathbf{q}(\mathbf{r})$  on  
 502 the network parameters  $W$ . We can redefine the loss,  $L^*$ , only leaving the terms that depends  
 503 on the trainable parameters  $W$ , and then apply a Monte Carlo method<sup>57</sup> to approximate the loss  
 504 from samples:

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

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

505

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

508 Therefore, by optimizing the overall function  $\mathbf{q}(\mathbf{r})$  to match the posterior distribution through  
 509 the use of cross-entropy loss, the network output  $f(\mathbf{r})$  can approximate the log of the likelihood  
 510 function  $\mathcal{L}(\theta)$  for each  $\mathbf{r}$  up to an unknown constant  $b(\mathbf{r})$ . Because we do not know the value of

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

As an example, consider a neural population with responses that follows a Poisson-like distribution (i.e. a version of the exponential distribution with linear sufficient statistics<sup>9,10</sup>). Learning a decoder for such population responses occurs as a special case of training a DNN-based likelihood decoder. For Poisson-like variability, the stimulus-conditioned distribution over  $\mathbf{r}$  is  $p(\mathbf{r}|\theta) = \phi(\mathbf{r})e^{\mathbf{h}^\top(\theta)\mathbf{r}}$ . The log likelihood function is then  $\log \mathbf{L} = \log \phi(\mathbf{r}) + \mathbf{H}^\top \mathbf{r}$ , where  $\mathbf{H}$  is a 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  $b(\mathbf{r}) = -\log \phi(\mathbf{r})$ . Hence, if we used a simple fully connected network, training the network is equivalent to fitting the kernel function  $\mathbf{h}(\theta)$  of the Poisson-like distribution.

In this work, we modeled the mapping  $f(\mathbf{r})$  as a DNN with two hidden layers<sup>17</sup>, consisting of two repeating blocks of a fully connected layer of size  $N_h$  followed by a rectified linear unit (ReLU)<sup>16</sup> and a drop-out layer<sup>58</sup> with dropout rate  $d_r$ , and a fully connected readout layer with no output nonlinearity (Fig. 3c). To encourage smoother likelihood functions, we added an  $L_2$  regularizer on  $\log \mathbf{L}$  filtered with a Laplacian filter of the form  $\mathbf{h} = [-0.25, 0.5, -0.25]$ . Therefore, the training loss included the term:

$$R = \gamma \sum_i \mathbf{u}_i^2 \quad (12)$$

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

We trained a separate instance of the network for each contrast-session, and referred to this class of DNN based likelihood decoder as the Full-Likelihood decoder to differentiate from alternative decoders described later.

During the training, each contrast-session was randomly split in proportions of 80% / 20% to yield the training set and the validation set, respectively. The stimulus orientation  $\theta$  was

binned into integers in the range  $[-45^\circ, 45^\circ]$ , and we excluded trials with orientations outside this range. This led to the exclusion of 157 out of 110,695 trials (0.14%) and 254 out of 192,631 trials (0.13%) for Monkey L and T data, respectively. The network was trained on the training set, starting with initial learning rate of  $\lambda_0$  and its performance on the validation set was monitored to perform early stopping<sup>59</sup>, and subsequently hyperparameter selection. For early stopping, we computed the mean squared error (MSE) between the maximum-a-posteriori (MAP) readout of the network output posterior  $\mathbf{q}$  and the stimulus orientation  $\theta$  on the validation set, and the training under a particular learning rate was terminated (early-stopped) if MSE failed to improve over 400 consecutive epochs, where each epoch is defined as one full pass through the training set. Upon early stopping, the parameter set that yielded the best validation set MSE during the course of the training was restored. The restored network was then trained again but with an updated learning rate  $\lambda_i = \frac{1}{3}\lambda_{i-1}$ , employing the same early stopping criteria. This procedure was repeated 4 times, therefore training the network under the 4 sequentially decreasing learning rate schedule of  $\lambda_0$ ,  $\frac{1}{3}\lambda_0$ ,  $\frac{1}{9}\lambda_0$  and  $\frac{1}{27}\lambda_0$ . Once the training was complete, the trained network was evaluated on the validation set to yield the final score, which served as the basis for our hyperparameter selections. The values of hyperparameters for the networks, including the size of the hidden layers  $N_h$ , the initial learning rate  $\lambda_0$ , the weight on the likelihood function smoothness regularizer  $\gamma$ , and the drop-out rate  $d_r$  during the training were selected by performing a random grid search over candidate values to find the combination that yielded the best validation set score for each contrast-session instance of the network (Supplementary Table 1). We observed that all possible values of hyperparameters were found among the best selected hyperparameter networks across all contrast-sessions and all types of networks trained.

Symbol	Description	Possible Values
$N_h$	number of hidden units per layer	$\{400, 600, 800, 1000\}$
$\lambda_0$	initial learning rate	$\{0.01, 0.03, 0.6\}$
$\gamma$	Laplacian L2 regularizer weight	$\{3, 30, 300\}$
$d_r$	dropout rate	$\{0.2, 0.5, 0.9\}$

Supplementary Table 1: Possible values of hyperparameters during model selection.

## 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<sup>14,18,19,21,22</sup>.

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

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

where  $\sigma_1 = 3^\circ$  and  $\sigma_2 = 15^\circ$ .

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

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

For the simulation, we have set  $A = 6$  and  $\sigma_{\text{sim}} = 21^\circ$ . We let the mean of the Gaussian tuning

571 curves for the 96 units to uniformly tile the stimulus orientation between  $-40^\circ$  and  $40^\circ$ . In other  
572 words,

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

573

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

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

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

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

578

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

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

582

583 for

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



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 this distribution, the variance of any unit's response scales linearly with its mean just as in the case of the Poisson distribution, but the population responses can be highly correlated depending on the choice of the correlation matrix  $C$ . For the simulation, we randomly generated a correlation matrix with the average units correlation of 0.227.

For each case of the distribution, we simulated population responses for the total of 1200 trials. Among these, 200 trials were set aside as the test set. We trained the DNN-based likelihood decoder on the remaining 1000 trials, splitting them further into 800 and 200 trials as the training and validation set, respectively. We followed the exact DNN training and hyperparameter selection procedure as described earlier.

For comparison, we also decoded the likelihood function from the population response  $\mathbf{r}_{\text{sim}}$  under the assumption of independent Poisson variability, regardless of the “true” distribution. Each unit's responses over the 1000 trials were fitted separately with a Gaussian tuning curve (Eq. 14). The parameters of the tuning curve  $A_i$ ,  $\mu_i$  and  $\sigma_{\text{sim}, i}$  were obtained by minimizing the least square difference between the Gaussian tuning curve and the observed  $i^{\text{th}}$  unit's responses  $(\theta, r_{\text{sim}, i})$  using `least_squares` function from Python SciPy optimization library.

The ground-truth likelihood function  $p(\mathbf{r}_{\text{sim}}|\theta)$  was computed for each simulated trial according to the definition of the distribution as found in Eq. 16 for the independent Poisson population or Eq. 18 for the mean scaled correlated Gaussian population.

We then assessed the quality of the decoded likelihood functions under the independent Poisson model  $\mathcal{L}_{\text{Pois}}(\theta)$  and under the DNN model  $\mathbf{L}_{\text{DNN}}$  by computing their Kullback-Leibler (KL) divergence<sup>56</sup> from the ground-truth likelihood function  $\mathcal{L}_{\text{gt}}(\theta)$ , giving rise to  $D_{\text{Pois}}$  and  $D_{\text{DNN}}$ , respectively. All continuous likelihood functions ( $\mathcal{L}_{\text{gt}}$  and  $\mathcal{L}_{\text{Pois}}$ ) were sampled at orientation  $\theta$  where  $\theta \in \mathbb{Z}$  and  $\theta \in [-45^\circ, 45^\circ]$ , giving rise to discretized likelihood functions  $\mathbf{L}_{\text{gt}}$  and

$\mathbf{L}_{\text{Poiss}}$  matching the dimensionality of the discretized likelihood function  $\mathbf{L}_{\text{DNN}}$  computed by the DNN. We then computed the KL divergence as:

$$D_{\text{Poiss}} = \sum_i \log \frac{\mathbf{L}_{\text{gt},i}}{\mathbf{L}_{\text{Poiss},i}} \mathbf{L}_{\text{gt},i} \quad (20)$$

and

$$D_{\text{DNN}} = \sum_i \log \frac{\mathbf{L}_{\text{gt},i}}{\mathbf{L}_{\text{DNN},i}} \mathbf{L}_{\text{gt},i} \quad (21)$$

We computed the KL divergence for both models across all 200 trials in the test set for both simulated population distributions (Supplementary Fig. 3). When the simulated population distribution was independent Poisson, then  $D_{\text{Poiss}} < D_{\text{DNN}}$  for all test set trials, indicating that  $\mathbf{L}_{\text{Poiss}}$  better approximated  $\mathbf{L}_{\text{gt}}$  overall than  $\mathbf{L}_{\text{DNN}}$ . However,  $\mathbf{L}_{\text{DNN}}$  still closely approximated  $\mathbf{L}_{\text{gt}}$ .

When the simulated population distribution was mean scaled correlated Gaussian,  $\mathbf{L}_{\text{DNN}}$  better approximated  $\mathbf{L}_{\text{gt}}$  than  $\mathbf{L}_{\text{Poiss}}$  on the majority of the trials. Furthermore,  $\mathbf{L}_{\text{Poiss}}$  provided qualitatively worse fit to the  $\mathbf{L}_{\text{gt}}$  for the simulated correlated Gaussian distribution compared to the fit of  $\mathbf{L}_{\text{DNN}}$  to  $\mathbf{L}_{\text{gt}}$  for the simulated independent Poisson distribution.

Overall, the simulation results suggest that (1) when the form of the underlying population distribution is known (such as in the case of independent Poisson population), more accurate likelihood functions can be decoded by directly using the knowledge of the population distribution than through the DNN-based likelihood decoder, but (2) when the form of the underlying distribution is unknown (such as in the case of the mean scaled correlated Gaussian distribution), then a DNN-based likelihood decoder can yield much more accurate likelihood functions than if one was to employ a wrong assumption about the underlying distribution in decoding likelihood functions, and (3) a DNN-based likelihood decoder can provide reasonable estimate

of the likelihood function across wide range of underlying distributions. Because the true underlying population distribution is hardly ever known to the experimenter, we believe that our DNN-based likelihood decoder stands as the most flexible method in decoding likelihood functions from the population responses to stimuli.

### **Fixed-Uncertainty likelihood decoder**

To test whether the trial-by-trial fluctuations in the shape of the likelihood function convey behaviorally relevant information, we developed the Fixed-Uncertainty likelihood decoder — a neural network based likelihood decoder that learns a fixed shape likelihood function whose location is shifted based on the input population response.

The Fixed-Uncertainty decoder network consisted of two parts: a learned fixed shape likelihood function  $L_0$  and a DNN that reads out a single scalar value  $\Delta_s$  corresponding to the shift that is applied to  $L_0$  (Supplementary Fig. 5) to yield the final likelihood function  $L$ . The DNN consisted of two repeating blocks of a fully connected layer followed by ReLU and a drop-out layer, and a final fully connected readout layer with no output nonlinearity, much like the DNN used for the Full-Likelihood decoder. The  $\log L_0$  was shifted by  $\Delta_s$  utilizing linear interpolation based grid-sampling<sup>60</sup> to shift the log-likelihood function in a manner that allows for the gradient of the loss to flow back to both the shift value  $\Delta_s$  (and therefore to the DNN parameters) as well as to the likelihood function shape  $L_0$ .

The output shifted log-likelihood function was then trained in an identical manner to the full-likelihood decoder described earlier, utilizing the same set of training paradigm with early stopping and regularizers, and explored the same range of hyperparameters.

## Likelihood functions based on Poisson-like and independent Poisson distributions

To serve as a comparison, for each contrast-session, we decoded likelihood functions from the population response assuming Poisson-like or independent Poisson distribution for  $p(\mathbf{r}|\theta)$  (Supplementary Fig. 2).

As was noted above, decoding likelihood function under the Poisson-like distribution is a special case of the Full-Likelihood decoder but using entirely linear DNN (i.e. no nonlinearity utilized in the network). Therefore, to decode likelihood functions under the assumption of the Poisson-like distribution, for each contrast-session, we trained a DNN with two hidden layers consisting of two repeating blocks of a fully connected layer followed by a drop-out layer<sup>58</sup> but with no nonlinear activation functions, and a fully connected readout layer with no output nonlinearity. The rest of the training and model selection procedure was identical to that of the Full-Likelihood or the Fixed-Uncertainty decoder described earlier.

To decode likelihood function under the independent Poisson distribution assumption, we first fitted tuning curves  $f_i(\theta)$  for each multiunit's responses to stimulus orientations  $\theta$  within a single contrast session. Tuning curves were computed using Gaussian process regression<sup>61</sup> with squared-exponential covariance function  $\text{cov}(f(\theta_1), f(\theta_2)) = \exp(-\frac{1}{2\sigma_L}(\theta_1 - \theta_2)^2)$  and a fixed observational noise  $\sigma_o$  using values of  $\sigma_L = 20$  and  $\sigma_o = 2$  selected based on the cross validation performance on multiunit's response prediction on a dataset not included elsewhere in the analysis. Once tuning curves were computed, the likelihood function over stimulus orientations was computed from the population response  $\mathbf{r}$  as follows:

$$\mathcal{L}(\theta) = \prod_i p(r_i|\theta) = \prod_i \frac{f_i(\theta)^{r_i} e^{-f_i(\theta)}}{r_i!} \quad (22)$$

## Mean and standard deviation of likelihood function

For uses in the subsequent analyses, we computed the mean and the standard deviation of the likelihood function by treating the likelihood function as an unnormalized probability distribution:

$$\mu_L = \frac{\int \theta \mathcal{L}(\theta) d\theta}{\int \mathcal{L}(\theta) d\theta} \quad (23)$$

$$\sigma_L = \sqrt{\frac{\int (\theta - \mu_L)^2 \mathcal{L}(\theta) d\theta}{\int \mathcal{L}(\theta) d\theta}} \quad (24)$$

We took the  $\mu_L$  and  $\sigma_L$  to be the point estimate of the stimulus orientation and the measure of the spread of the likelihood function, respectively, used in all subsequent analyses. Although not presented here, we performed the following analyses with other point estimates of the stimulus orientation such as the orientation at the maximum of the likelihood function and the median of the likelihood functions, and observed that models with mean of the likelihood function as the point estimate performed the best.

## Attribution analysis

To assess whether the same members of the population simulatenously encode the best point estimate (i.e. in the form of the mean of the likelihood function  $\mu_L$ ) and uncertainty (i.e. in the form of the width of the likelihood function  $\sigma_L$ ), we computed the attribution of each multiunit input of the trained Full-Likelihood decoder to the mean of the likelihood  $\mu_L$  and the standard deviation of the likelihood function  $\sigma_L$  giving rise to the attribution  $A_\mu, A_\sigma \in \mathbb{R}^m$ , respectively, where  $m$  is the number of multiunits in the input to the network. Among numerous methods of computing attribution<sup>33–35,62</sup>, we have selected three popular gradient based attribution methods<sup>33</sup>: saliency maps<sup>34</sup>, gradient  $\times$  input<sup>62</sup>, and DeepLift<sup>35</sup> and compared their results.

Given a collection of input population responses and computed likelihood functions  $\{\mathbf{r}^{(i)}, \mathbf{L}^{(i)}\}$ , where the superscript denotes the  $i^{\text{th}}$  trial in the contrast session, we compute the mean and the standard deviation of the likelihood function according to Eq. 23 and Eq. 24, respectively, giving rise to  $\mu_L^{(i)}$  and  $\sigma_L^{(i)}$ . Given a target feature  $S \in \{\mu_L, \sigma_L\}$  that can be computed from the input units  $\mathbf{r}$  through a differentiable function, we compute the attribution of the input units to the target  $S$  for each trial according to each attribution method, yielding  $\mathbf{a}_{S,\text{method}}^{(k)}$ , where  $\mathbf{a} \in \mathbf{R}^m$ . The sign of the attribution indicates whether increasing the unit tends to increase or decrease the target feature. Since we are interested more in how much each unit contribute to the target feature rather than in which direction, we take the absolute value of per trial attribution and compute the average across all trials to yield the final attribution of the input units:

$$\mathbf{A}_{S,\text{method}} = \sum_k |a_{S,\text{method}}^{(k)}| \quad (25)$$

For the saliency maps based method<sup>34</sup>, the attribution is computed as the partial derivative of the feature  $S$  with respect to the input units  $\mathbf{r}$ :

$$\mathbf{a}_{S,\text{Saliency}} = \frac{\partial S}{\partial \mathbf{r}} \quad (26)$$

which can be computed rather straightforwardly on a DNN implemented using any of the modern neural network libraries equipped with automatic gradient computation.

For Gradient  $\times$  Input (GI) method, the attribution is computed as the gradient of the feature with respect to the input (as in saliency maps) multiplied with the input  $\mathbf{r}$ :

$$\mathbf{a}_{S,\text{GI}} = \frac{\partial S}{\partial \mathbf{r}} \odot \mathbf{r} \quad (27)$$

Finally, we computed DeepLIFT attribution by using modified gradient computation for ReLU units in the network defined as:

$$\frac{\partial^m \text{ReLU}(x)}{\partial x} = \frac{\text{ReLU}(x) - \text{ReLU}(x_0)}{x - x_0} \quad (28)$$

where  $x_0$  represents the input into the ReLU nonlinearity when a reference input  $\mathbf{r}_0$  was used as the input into the network. Here, we have defined the reference network input to be the average population response across all trials (refer to Ref<sup>33,35</sup> for details).

Using the above modified gradient computation for ReLU nonlinearity in the backpropagation to compute the partial derivative of the target feature with respect to the input units yield the modified partial derivative  $\frac{\partial^m S}{\partial \mathbf{r}}$  which is finally used to compute the DeepLIFT (DL) attribution as:

$$\mathbf{a}_{S,DL} = \frac{\partial^m S}{\partial \mathbf{r}} \odot (\mathbf{r} - \mathbf{r}_0) \quad (29)$$

For each contrast session and each attribution method, we computed the attribution of the input units to both  $\mu_L$  and  $\sigma_L$ , yielding vectors  $\mathbf{A}_\mu$  and  $\mathbf{A}_\sigma$ , and we computed the Pearson correlation coefficient between the two scores across the units (Fig. 6).

## Decision-making models

Given the hypothesized representation of the stimulus and its uncertainty in the form of the likelihood function  $\mathcal{L}(\theta) \equiv p(\mathbf{r}|\theta)$ , the monkey's trial-by-trial decisions were modeled based on the assumption that the monkey computes the posterior probability over the two classes  $C = 1$  and  $C = 2$ , and utilizes this information in making decisions—that is, in accordance to a model of a Bayesian decision maker. The orientation distributions for the two 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  $\sigma_1 = 3^\circ$  and  $\sigma_2 = 15^\circ$  where  $\mathcal{N}(\theta; \mu, \sigma^2)$  denotes a Gaussian distribution over  $\theta$  with mean  $\mu$  and variance  $\sigma^2$ . The posterior ratio  $\rho$  for the two classes is:

$$\rho = \frac{p(C = 2|\mathbf{r})}{p(C = 1|\mathbf{r})} \quad (30)$$

$$= \frac{p(C = 2) \int p(\mathbf{r}|\theta)p(\theta|C = 2) d\theta}{p(C = 1) \int p(\mathbf{r}|\theta)p(\theta|C = 1) d\theta} \quad (31)$$

$$= \frac{p(C = 2) \int \mathcal{L}(\theta)\mathcal{N}(\theta; \mu, \sigma_2^2) d\theta}{p(C = 1) \int \mathcal{L}(\theta)\mathcal{N}(\theta; \mu, \sigma_1^2) d\theta} \quad (32)$$

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

$$\hat{C} = \underset{C}{\operatorname{argmax}} p(C|\mathbf{r}) \quad (33)$$

where  $\hat{C}$  is the subject's decision. However, according to the posterior probability matching strategy<sup>63,64</sup>, the decision of subjects on certain tasks are better modeled as sampling from the posterior probability:

$$p(\hat{C}) = p(C = \hat{C}|\mathbf{r}) \quad (34)$$

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

$$\frac{p(\hat{C} = 2)}{p(\hat{C} = 1)} = \left( \frac{p(C = 2|\mathbf{r})}{p(C = 1|\mathbf{r})} \right)^\alpha = \rho^\alpha \quad (35)$$

where  $\alpha \in \mathbb{R}^+$ . When  $\alpha = 1$ , the decision-making strategy corresponds to the posterior probability matching while  $\alpha = \infty$  corresponds to the MAP strategy<sup>64</sup>. We fitted the value of  $\alpha$



for each contrast-session during the model fitting to capture any variation of the strategy. Furthermore, we incorporated a lapse rate  $\lambda$ , a fraction of trials on which the subject does not pay attention and makes a random decision. Hence, the final probability that the subject selects the class  $C = 1$  was modeled as:

$$p(\hat{C} = 1) = (1 - \lambda) \frac{1}{1 + \rho^\alpha} + 0.5\lambda \quad (36)$$

$$= (1 - \lambda) \left[ 1 + \left( \frac{p(C = 2) \int \mathcal{L}(\theta) \mathcal{N}(\theta; \mu, \sigma_2^2) d\theta}{p(C = 1) \int \mathcal{L}(\theta) \mathcal{N}(\theta; \mu, \sigma_1^2) d\theta} \right)^\alpha \right]^{-1} + 0.5\lambda \quad (37)$$

$$= (1 - \lambda) \left[ 1 + \left( \frac{(1 - p(C = 1)) \int \mathcal{L}(\theta) \mathcal{N}(\theta; \mu, \sigma_2^2) d\theta}{p(C = 1) \int \mathcal{L}(\theta) \mathcal{N}(\theta; \mu, \sigma_1^2) d\theta} \right)^\alpha \right]^{-1} + 0.5\lambda \quad (38)$$

For each contrast-session, we fitted the above Bayesian decision model to the monkey's decisions by fitting the four parameters:  $\mu$ ,  $p(C = 1)$ ,  $\alpha$ , and  $\lambda$ . Fitting  $\mu$  (the center of stimulus orientation distributions) and  $p(C = 1)$  (prior over class) allowed us to capture the bias in the stimulation distribution that the subject may have acquired erroneously during the training, and fitting  $\alpha$  and  $\lambda$  allowed for the model to match the decision-making strategy employed by the subject.

Utilizing the likelihood function  $\mathcal{L}(\theta)$  decoded from the V1 population response via the Full-Likelihood decoder network in Eq. 38 gave rise to the Full-Likelihood Model that made use of all information including the trial-by-trial uncertainty information as captured by the trial-by-trial fluctuations in the shape of the likelihood function. Alternatively, utilizing the likelihood function decoded by the trained Fixed-Uncertainty decoder gave rise to the Fixed-Uncertainty Model. The Fixed-Uncertainty Model effectively ignores all trial-by-trial fluctuations in the uncertainty that would be captured by the fluctuations in the shape of the likelihood function, but captures the trial-by-trial point estimate of the stimulus orientation  $\hat{\theta}$  by shifting the learned fixed shape likelihood function over orientation. For each contrast-session, different fixed likelihood

shape was learned, allowing the overt measure of uncertainty such as contrast to modulate the expected level of uncertainty.

For comparison, we have also tested the performance of the trial-by-trial decision prediction utilizing likelihood functions decoded based on Poisson-like or independent Poisson population distribution assumptions, giving rise to the Poisson-like Model and the Independent Poisson Model for predicting trial-by-trial decisions, respectively.

## 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,  $B_1, B_2, \dots, B_i, \dots, B_{10}$  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:

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \left( -\log \prod_i p(\hat{C} = \hat{C}_i | M, \Theta) \right) \quad (39)$$

$$= \underset{\Theta}{\operatorname{argmin}} \left( -\sum_i \log p(\hat{C} = \hat{C}_i | M, \Theta) \right) \quad (40)$$

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. 38 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)<sup>65</sup>. 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. Specifically, we binned trials to the nearest orientation degree such that each bin was centered at an integer degree (i.e. bin center  $\in \mathbb{Z}$ ) with the bin width of  $1^\circ$ . We then shuffled the likelihood functions among trials in the same orientation bin. This effectively removed the stimulus orientation conditioned correlation between the likelihood function and the subject's classification  $\hat{C}$ , while preserving the expected likelihood function for each stimulus orientation.

However, we were specifically interested in decoupling the uncertainty information contained in the shape of the likelihood function from the decision while minimally disrupting the trial-by-trial correlation between the point estimate of the stimulus orientation and the subject's classification decision. To achieve this, for each trial, the newly assigned likelihood function

was shifted such that the mean of the normalized likelihood function,  $\mu_L$  (Eq. 23), remained the same for each trial before and after the likelihood shuffling (Fig. 5c). This allowed us to specifically assess the effect of distorting the shape of the likelihood function conditioned on both the (binned) stimulus orientation and the point estimate of the stimulus orientation (i.e.  $\mu_L$ ) (Fig. 5c). To ensure that both models can take the full advantage of any information that remains in the shuffled likelihood functions, we trained both the Full-Likelihood Model and the Fixed-Uncertainty Model from scratch on the shuffled data. Aside from the difference in the dataset, we followed the exact procedure used when training on the original (unshuffled) data, evaluating the performance through cross-validation on the test sets.

## Classification simulation

We computed the expected effect size of the model fit difference between the Full-Likelihood Model and the Fixed-Uncertainty Model by generating simulated data using the trained Full-Likelihood Model as the ground truth. Specifically, for each trial for each contrast-session, we computed the probability of responding  $\hat{C} = 1$  from Eq. 38, utilizing the full decoded likelihood function  $\mathcal{L}(\theta)$  for the given trial, and sampled a classification decision from that probability. This procedure yielded simulated data where all monkeys' decisions were replaced by decisions made by the trained Full-Likelihood Models. We repeated this procedure 5 times, thereby producing 5 sets of simulated data. For each set of simulated data, we trained the two decision-making models (Full-Likelihood Model and Fixed-Uncertainty Model) on each contrast-session with 10-fold cross-validation, and reported the aggregated test set scores as was done for the original data.

## Code availability

Code used for modeling and training the deep neural networks as well as for figure generation will be made available for view and download at [https://github.com/eywalker/v1\\_likelihood](https://github.com/eywalker/v1_likelihood). All other code used for analysis including data selection and decision model fitting will be placed at [https://github.com/eywalker/v1\\_project](https://github.com/eywalker/v1_project). Finally, code used for electrophysiology data processing can already be found in the Tolias lab GitHub organization <https://github.com/atlab>.

## Data availability

All figures except for Figure 1 and Supplementary Figure 4 were generated from raw data or processed data. The data generated and/or analyzed during the current study are available from the corresponding author upon reasonable request. No publicly available data was used in this study.

## Statistics

All statistical tests used were two-tailed paired two-sample t-test, unless specified otherwise. Wherever reported, data are means and error bars indicate standard error of the means computed as  $\frac{\sigma}{\sqrt{n}}$  where  $\sigma$  is the standard deviation and  $n$  is the size of the sample within the bin, unless specified otherwise. Exact p values less than 0.001 were reported as  $p < 0.001$ . When appropriate, p values were corrected for multiple comparisons and the corrected p value was reported.

## Acknowledgments

The research was supported by National Science Foundation Grant IIS-1132009 (to W.J.M. and A.S.T.), DP1 EY023176 Pioneer Grant (to A.S.T.), F30 EY025510 (to E.Y.W.) and R01

EY026927 (to A.S.T and W.J.M.). We thank Fabian Sinz for helpful discussion and suggestions on the deep neural network fitting to likelihood functions. We also thank Tori Shin for assistance in monkeys behavioral training and experimental data collection.

## Author 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.

## 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. 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).

- 886 16. LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436–444 (2015).
- 887 17. Goodfellow, Ian, Bengio, Yoshua, Courville, A. Deep Learning. *MIT Press* (2016).
- 888     arXiv:1312.6184v5.
- 889 18. Seung, H. S. & Sompolinsky, H. Simple models for reading neuronal population codes.
- 890     *Proc.Natl.Acad.Sci.* **90**, 10749–10753 (1993). NIHMS150003.
- 891 19. Sanger, T. D. Probability density estimation for the interpretation of neural population
- 892     codes. *Journal of neurophysiology* **76**, 2790–3 (1996).
- 893 20. Zemel, R. S., Dayan, P. & Pouget, A. Probabilistic interpretation of population codes.
- 894     *Neural Comp.* **10**, 403–430 (1998).
- 895 21. Jazayeri, M. & Movshon, J. A. Optimal representation of sensory information by neu-
- 896     ral populations. *Nature Neuroscience* **9**, 690–696 (2006). 10.1021/nl3012853 |
- 897     NanoLett.2012, 12, 36023608.
- 898 22. Fetsch, C. R., Pouget, A., Deangelis, G. C. & Angelaki, D. E. Neural correlates of
- 899     reliability-based cue weighting during multisensory integration. *Nature Neuroscience* **15**,
- 900     146–154 (2012). NIHMS150003.
- 901 23. Averbeck, B. B. & Lee, D. Effects of Noise Correlations on Information Encoding and
- 902     Decoding. *J Neurophysiol* **95**, 3633–3644 (2006).
- 903 24. Ecker, A. S. *et al.* Decorrelated neuronal firing in cortical microcircuits. *Science* **327**,
- 904     584–587 (2010).
- 905 25. Ecker, A. S., Berens, P., Tolias, A. S. & Bethge, M. The Effect of Noise Correlations
- 906     in Populations of Diversely Tuned Neurons. *Journal of Neuroscience* **31**, 14272–14283
- 907     (2011). NIHMS150003.



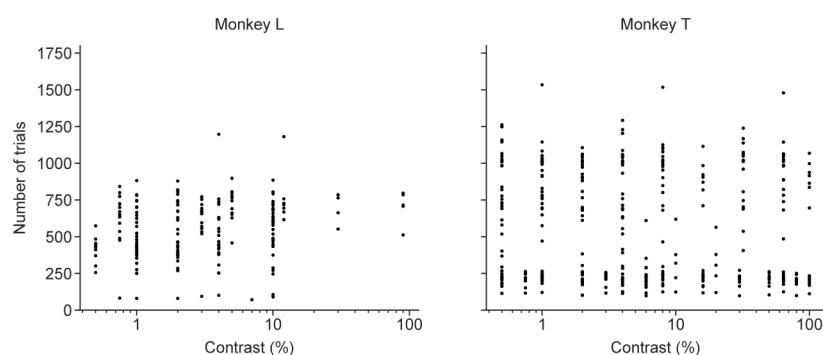
- 908 26. Ecker, A. S. *et al.* State dependence of noise correlations in macaque primary visual cortex.  
909 *Neuron* **82**, 235–248 (2014).
- 910 27. van Bergen, R. S. & Jehee, J. F. Modeling correlated noise is necessary to decode uncer-  
911 tainty. *NeuroImage* (2017). 1708.04860.
- 912 28. Denfield, G. H., Ecker, A. S., Shinn, T. J., Bethge, M. & Tolias, A. S. Attentional fluctua-  
913 tions induce shared variability in macaque primary visual cortex. *Nature Communications*  
914 **9**, 2654 (2018).
- 915 29. Ma, W. J. Signal detection theory, uncertainty, and Poisson-like population codes. *Vision*  
916 *Research* **50**, 2308–2319 (2010).
- 917 30. Van Bergen, R. S., Ji Ma, W., Pratte, M. S. & Jehee, J. F. Sensory uncertainty decoded from  
918 visual cortex predicts behavior. *Nature Neuroscience* **18**, 1728–1730 (2015). 15334406.
- 919 31. Tolhurst, D. J., Movshon, J. A. & Dean, A. F. The statistical reliability of signals in single  
920 neurons in cat and monkey visual cortex. *Vision Research* **23**, 775–785 (1983).
- 921 32. Shadlen, M. N. & Newsome, W. T. The variable discharge of cortical neurons: implications  
922 for connectivity, computation, and information coding. *The Journal of neuroscience : the*  
923 *official journal of the Society for Neuroscience* **18**, 3870–96 (1998).
- 924 33. Ancona, M., Ceolini, E., Öztireli, C. & Gross, M. A unified view of gradient-based attri-  
925 bution methods for deep neural networks (2017).
- 926 34. Simonyan, K., Vedaldi, A. & Zisserman, A. Deep inside convolutional networks: Visualis-  
927 ing image classification models and saliency maps. *CoRR* **abs/1312.6034** (2013).
- 928 35. Shrikumar, A., Greenside, P. & Kundaje, A. Learning important features through propa-  
929 gating activation differences. In Precup, D. & Teh, Y. W. (eds.) *Proceedings of the 34th In-*

- ternational Conference on Machine Learning, vol. 70 of *Proceedings of Machine Learning Research*, 3145–3153 (PMLR, International Convention Centre, Sydney, Australia, 2017).
36. Campbell, F. & Kulikowski, J. The visual evoked potential as a function of contrast of a grating pattern. *The Journal of physiology* **222**, 345–356 (1972).
37. Britten, K. H., Newsome, W. T., Shadlen, M. N., Celebrini, S. & Movshon, J. A. A relationship between behavioral choice and the visual responses of neurons in macaque mt. *Visual neuroscience* **13**, 87–100 (1996).
38. Angelaki, D. E., Humphreys, G. & DeAngelis, G. C. Multisensory Integration. *Journal Of The Theoretical Humanities* **19**, 452–458 (2009).
39. 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.
40. Beck, J. M., Latham, P. E. & Pouget, A. Marginalization in Neural Circuits with Divisive Normalization. *Journal of Neuroscience* **31**, 15310–15319 (2011). NIHMS150003.
41. Ma, W. J. & Rahmati, M. Towards a Neural Implementation of Causal Inference in Cue Combination. *Multisensory Research* **26**, 159–176 (2013).
42. Orhan, A. E. & Ma, W. J. Efficient probabilistic inference in generic neural networks trained with non-probabilistic feedback. *Nature Communications* **8**, 138 (2017).
43. Cumming, B. G. & Nienborg, H. Feedforward and feedback sources of choice probability in neural population responses. *Current opinion in neurobiology* **37**, 126–132 (2016).
44. Bondy, A. G., Haefner, R. M. & Cumming, B. G. Feedback determines the structure of correlated variability in primary visual cortex. *Nature neuroscience* **21**, 598 (2018).

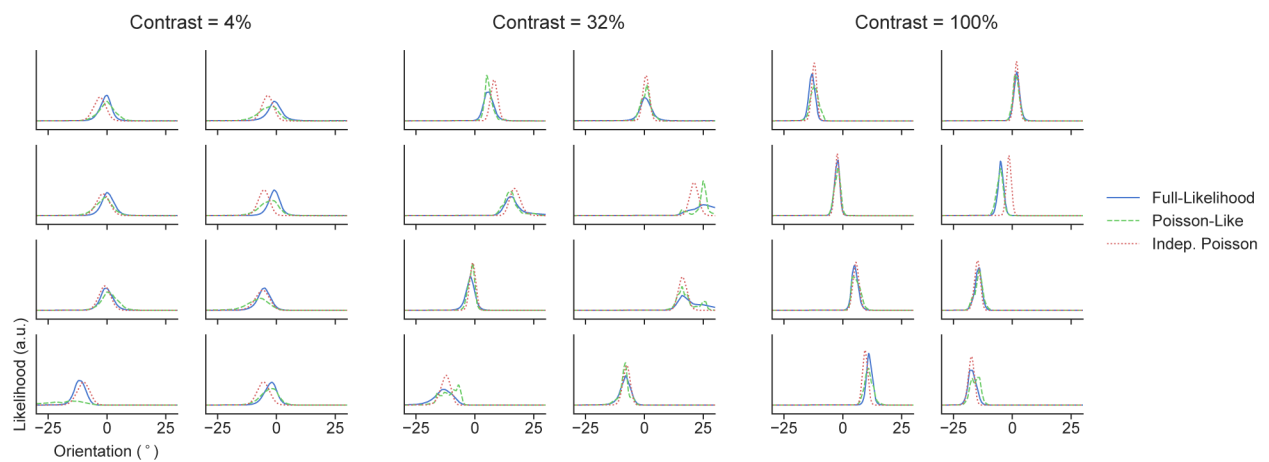
45. Geisler, W. S. Contributions of ideal observer theory to vision research (2011).  
NIHMS150003.
46. Körding, K. Decision theory: What "should" the nervous system do? (2007).
47. Maloney, L. T. & Mamassian, P. Bayesian decision theory as a model of human visual perception: Testing Bayesian transfer. *Visual Neuroscience* **26**, 147–155 (2009).
48. Ma, W. J. Organizing probabilistic models of perception. *Trends in Cognitive Sciences* **16**, 511–518 (2012).
49. Brainard, D. H. The Psychophysics Toolbox. *Spatial vision* **10**, 433–6 (1997).
50. Tolias, A. S. *et al.* Recording Chronically From the Same Neurons in Awake, Behaving Primates. *Journal of Neurophysiology* **98**, 3780–3790 (2007).
51. Subramaniam, M., Ecker, A. S., Berens, P. & Tolias, A. S. Macaque Monkeys Perceive the Flash Lag Illusion. *PLoS ONE* **8**, e58788 (2013).
52. Quiroga, R. Q., Nadasdy, Z. & Ben-Shaul, Y. Unsupervised Spike Detection and Sorting with Wavelets and Superparamagnetic Clustering. *Neural Computation* **16**, 1661–1687 (2004).
53. Kohn, A. & Movshon, J. A. Adaptation changes the direction tuning of macaque MT neurons. *Nature Neuroscience* **7**, 764–772 (2004).
54. Richard, M. D. & Lippmann, R. P. Neural Network Classifiers Estimate Bayesian a posteriori Probabilities. *Neural Computation* **3**, 461–483 (1991).
55. 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).

- 974 56. Kullback, S. & Leibler, R. A. On Information and Sufficiency. *The Annals of Mathematical*  
975 *Statistics* **22**, 79–86 (1951). 1511.00860.
- 976 57. MacKay, D. J. C. *Information Theory , Inference , and Learning Algorithms*, vol. 22 (Cam-  
977 bridge University Press, Cambridge, UK, 2003). arXiv:1011.1669v3.
- 978 58. Srivastava, N., Hinton, G., Krizhevsky, A. & Salakhutdinov, R. Dropout: A Simple Way  
979 to Prevent Neural Networks from Overfitting. *Journal of Machine Learning Research* **15**,  
980 1929–1958 (2014).
- 981 59. Prechelt, L. Early stopping-but when? In *Neural Networks: Tricks of the Trade, This Book*  
982 *is an Outgrowth of a 1996 NIPS Workshop*, 55–69 (Springer-Verlag, London, UK, UK,  
983 1998).
- 984 60. Jaderberg, M., Simonyan, K., Zisserman, A. *et al.* Spatial transformer networks. In *Ad-*  
985 *vances in neural information processing systems*, 2017–2025 (2015).
- 986 61. Rasmussen, C. E. Gaussian processes in machine learning 63–71 (2003).
- 987 62. Shrikumar, A., Greenside, P., Shcherbina, A. & Kundaje, A. Not just a black box: Learn-  
988 ing important features through propagating activation differences. *CoRR* **abs/1605.01713**  
989 (2016). URL <http://arxiv.org/abs/1605.01713>. 1605.01713.
- 990 63. Mamassian, P. & Landy, M. S. Observer biases in the 3D interpretation of line drawings.  
991 *Vision Research* **38**, 2817–2832 (1998).
- 992 64. Acerbi, L., Vijayakumar, S. & Wolpert, D. M. On the Origins of Suboptimality in Human  
993 Probabilistic Inference. *PLoS Computational Biology* **10**, e1003661 (2014).

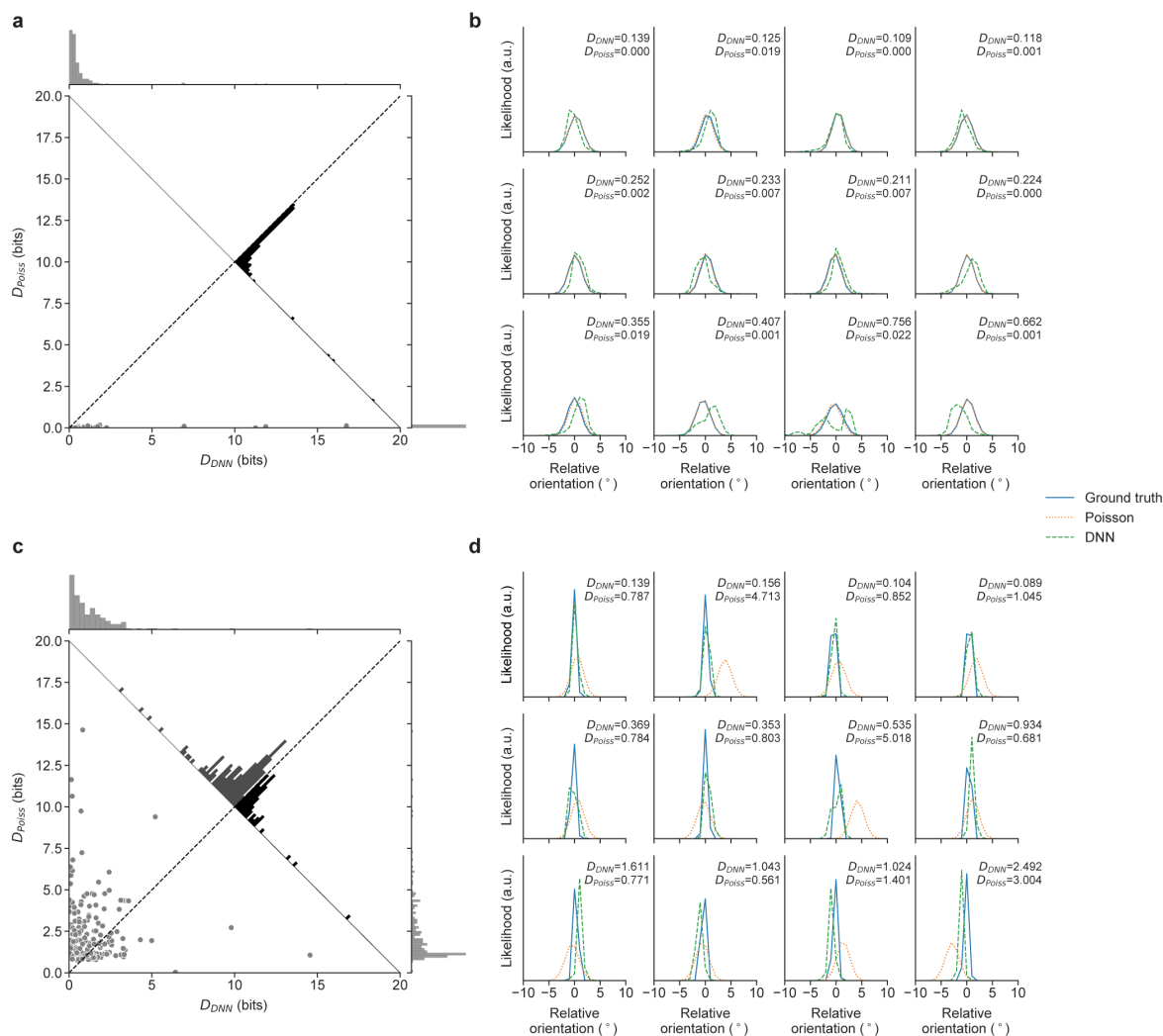
- 994 65. Acerbi, L. & Ma, W. J. Practical Bayesian Optimization for Model Fitting with Bayesian  
 995 Adaptive Direct Search. In *Advances in Neural Information Processing Systems 30*, 1836–  
 996 1846 (2017). 1705.04405.



Supplementary Figure 1: Number of trials per contrast-session. Each point corresponds to a single contrast-session, depicting the number of trials performed at the particular contrast.

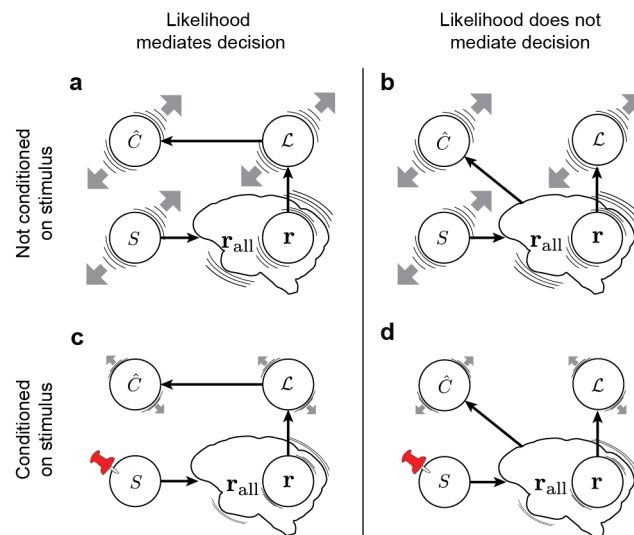


Supplementary Figure 2: Example decoded likelihood functions. Example decoded likelihood functions under Full-Likelihood, Poisson-like and Independent-Poisson based decoders are shown for randomly selected trials from three distinct contrast-sessions from Monkey T.

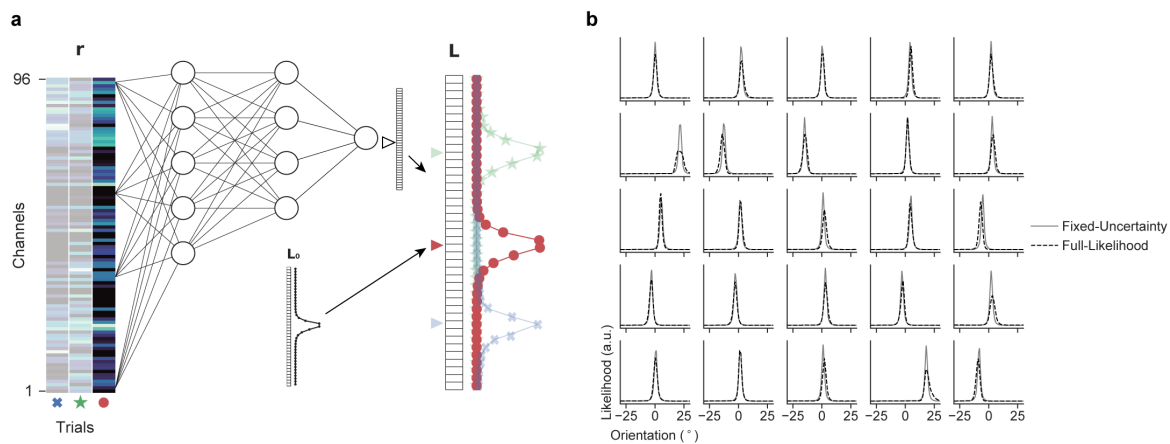


Supplementary Figure 3: Performance of the likelihood functions decoded by DNN-based decoders. **a-b**, Results on independent Poisson population responses. **a**, KL divergence between the ground truth likelihood function and likelihood function decoded with: a trained DNN  $D_{\text{DNN}}$  vs. independent 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.

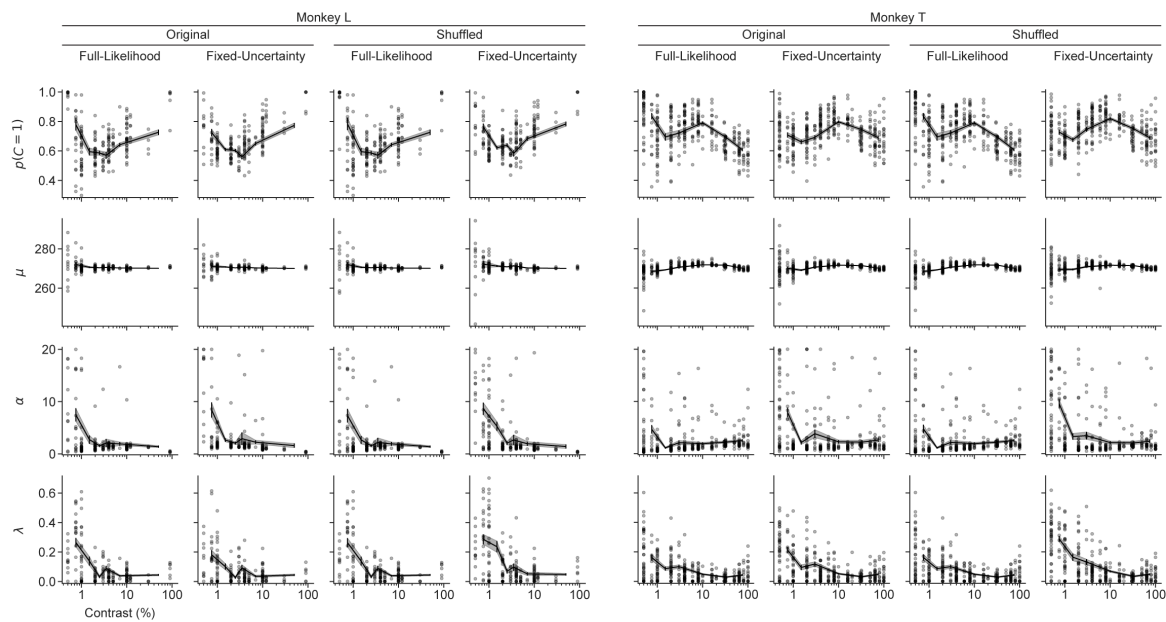




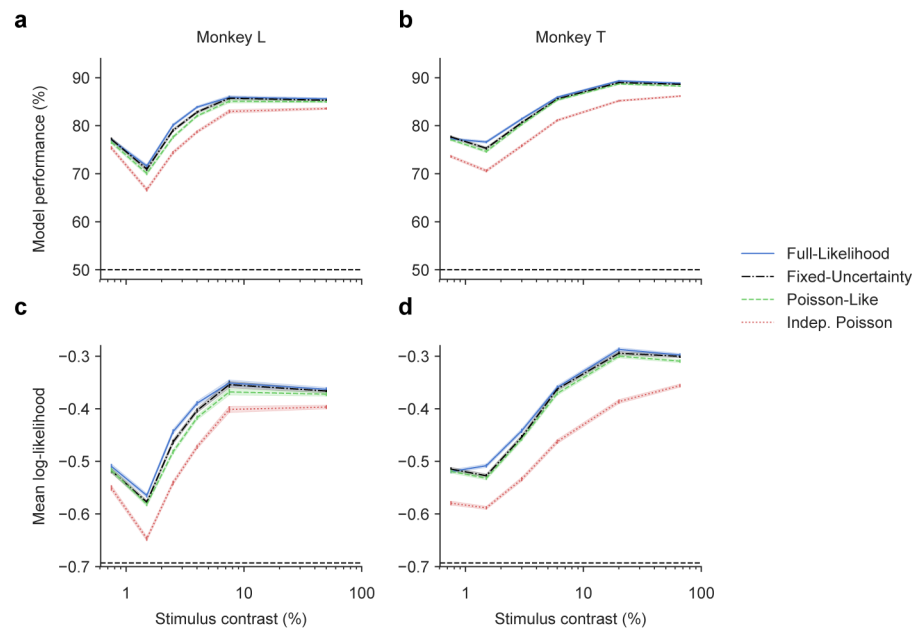
Supplementary Figure 4: Alternative relationships between the likelihood function and the decision. Possible relationships between variables in the model are indicated by black arrows. We consider two scenarios: **a, c** the likelihood function  $\mathcal{L}$  mediates the decision  $\hat{C}$ , **b, d** 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. **a, b**, When not conditioning on the stimulus  $s$ , the stimulus can drive correlation among all variables, making it difficult to distinguish the two scenarios. **c, d**, 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. The variable  $r$  represents the recorded cortical population and  $r_{\text{all}}$  represents responses of all recorded and unrecorded neurons.



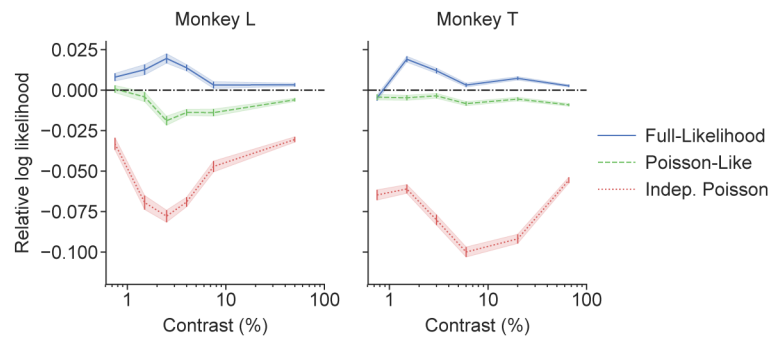
Supplementary Figure 5: Fixed-Uncertainty decoder. **a**, A schematic of a DNN for the Fixed-Uncertainty decoder mapping  $r$  to the decoded likelihood function  $L$ . For each contrast-session, the Fixed-Uncertainty decoder learns a single fixed-shape likelihood function  $L_0$  and a network that shifts  $L_0$  based on the population response. Therefore, all resulting likelihood functions share the same shape (uncertainty) but differ in the center location from trial to trail. **b**, Example decoded likelihood functions from randomly selected trials from a single contrast session for both the Fixed-Uncertainty decoder and Full-Likelihood decoder.



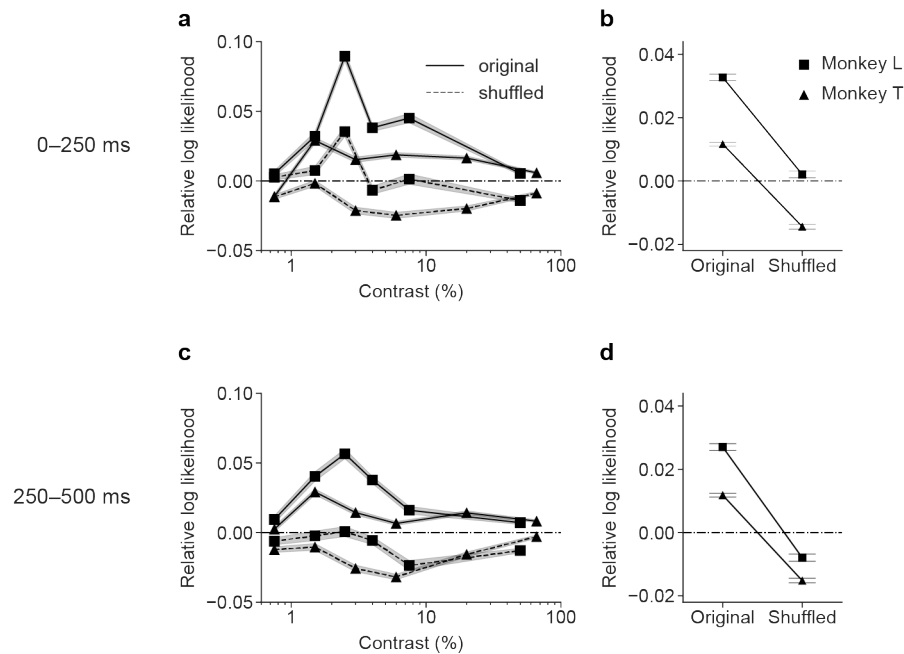
Supplementary Figure 6: Fitted Bayesian decision model parameters. Each point corresponds to a single contrast-session, depicting the average fitted parameter value across 10 cross-validation training sets plotted against the contrast of the contrast-session. The solid line and error bars/shaded area depicts the mean and the standard error of the mean of the parameter value for binned contrast values, respectively.



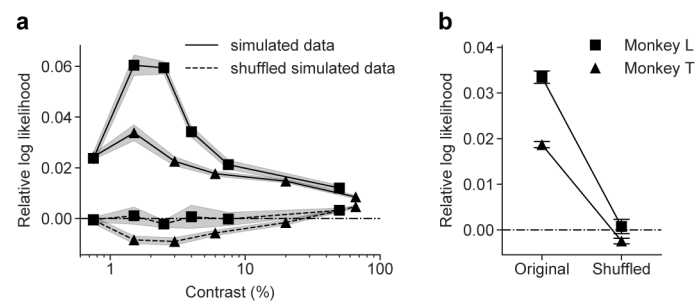
Supplementary Figure 7: Model performance on decision predictions. **a-b**, Model performance measured in proportions of trials correctly predicted by the model as a function of contrast for four decision models based on different likelihood decoders. On each trial, the class decision that would maximize the posterior  $p(\hat{C}|\mathbf{r})$  was chosen to yield a concrete classification prediction. **c-d**, Same as in **a-b** but with performance measured as the trial-averaged log likelihood of the model. For **a-b** and **c-d**, dashed lines indicate the performance at chance (50% and  $\ln(0.5)$ , respectively).



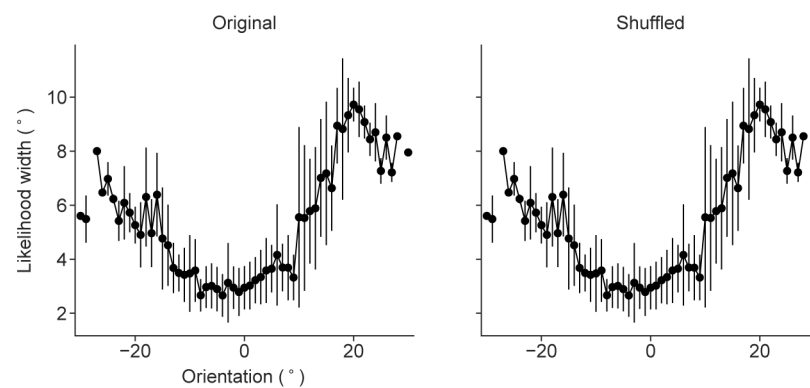
Supplementary Figure 8: Performance of Poisson-like and Independent Poisson Models. For each monkey, the average trial-by-trial performance of the Full-Likelihood, Poisson-like and Independent Poisson Models are shown relative to the Fixed-Uncertainty Model across contrasts, measured as the average trial difference in the log likelihood.



Supplementary Figure 9: Model performance based on population responses to different stimulus windows. **a, c**, 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 models were trained and evaluated on the population response to (a) the first half (0–250 ms) or (c) the second half (250–500 ms) of the stimulus presentation. 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, d**, Relative model performance summarized across all contrasts based on models trained as described in (a, c). Performance on the original and the shuffled data is shown individually for both monkeys. The difference between the Full-Likelihood and Fixed-Uncertainty Models was significant with  $p < 0.001$  for both stimulus windows, and on both the original and the shuffled data for both monkeys, except for the shuffled dataset on 0–250ms for Monkey L, for which there was no significant difference between the two models ( $p = 0.17$ ). The difference between the Full-Likelihood Model on the original and the shuffled data was significant ( $p < 0.001$  for both monkeys for both stimulus windows). For **a-d**, all data points are means, and error bar/shaded area indicate standard error of the means.



Supplementary Figure 10: 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 and shuffled 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.



Supplementary Figure 11: Dependence of average likelihood width on the stimulus orientation. The dependence of the width of the likelihood function  $\sigma_L$  on the stimulus orientation is depicted for an example contrast-session (Monkey T, 8% contrast) on the original and the shuffled data. The shuffling procedure preserves the relationship between the average likelihood width and the stimulus orientation as desired.