

A neural basis of probabilistic computation in visual cortex

Edgar Y. Walker^{1,2,6*}, R. James Cotton^{1,2,5,6}, Wei Ji Ma^{3,7*} and Andreas S. Tolias^{1,2,4,7*}

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 tested this hypothesis by simultaneously recording population activity from 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 Laplace¹ and von Helmholtz², has been widely used to explain human and animal perception^{3–6}. At its core lies the assumption that the brain maintains a statistical model of the world and when confronted with incomplete and imperfect information, it makes inferences by computing probability distributions over task-relevant world state variables (for example, 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 (for example, Alais and Burr⁷ and Ernst and Banks⁸). 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 (for example, 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^{9,10}. Specifically, according to the probabilistic population coding (PPC) hypothesis^{9,10}, inference in the brain is performed by inverting a generative model of neural population activity. Under this coding scheme, neural populations in V1, for example, that encode stimulus orientation also encode the associated uncertainty in the form of the sensory likelihood function, the probability of observing a given pattern of neural activity across hypothesized stimulus values^{9,11,12}. The form of the likelihood function is related to the probability distribution describing neural variability (noise) for a given stimulus. A sensory likelihood function is often unimodal^{13,14}, 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 (for example, 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 animal¹⁵. To decode the trial-by-trial likelihood functions from the V1 population responses, we developed a technique based on deep learning^{16,17}. 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 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 and is consistent with the theory of PPC.

Results

Behavior. Two Rhesus macaques (*Macaca 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

¹Center for Neuroscience and Artificial Intelligence, Baylor College of Medicine, Houston, TX, USA. ²Department of Neuroscience, Baylor College of Medicine, Houston, TX, USA. ³Center for Neural Science, Department of Psychology, New York University, New York, NY, USA. ⁴Department of Electrical and Computer Engineering, Rice University, Houston, TX, USA. ⁵Present address: Shirley Ryan AbilityLab, Chicago, IL, USA. ⁶These authors contributed equally: Edgar Y. Walker, R. James Cotton. ⁷These authors jointly supervised this work: Wei Ji Ma, Andreas S. Tolias. *e-mail: eywalker@bcm.edu; weijima@nyu.edu; astolias@bcm.edu

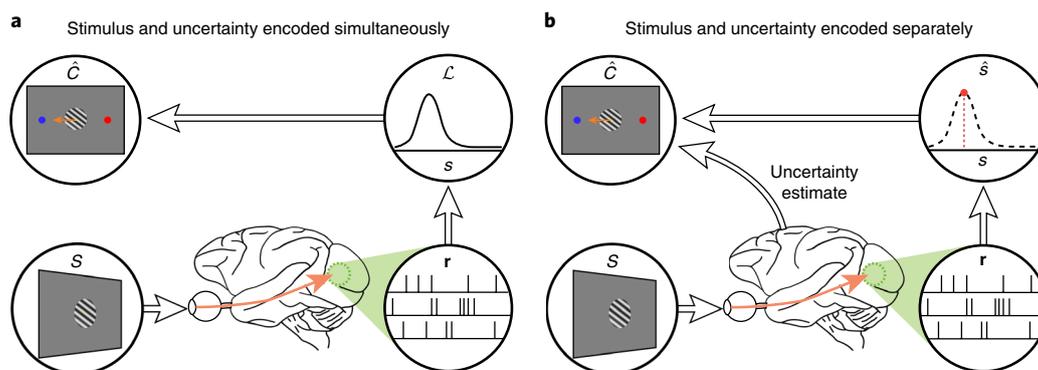


Fig. 1 | Alternative models of uncertainty information encoding. **a**, The recorded cortical population, r , responding to sensory stimulus, s , encodes stimulus estimate and uncertainty simultaneously in the form of likelihood function L , which is subsequently used in making a decision \hat{C} as the subject performs a visual classification task. **b**, The recorded cortical population encodes only a point estimate of the stimulus, \hat{s} , whereas 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} .

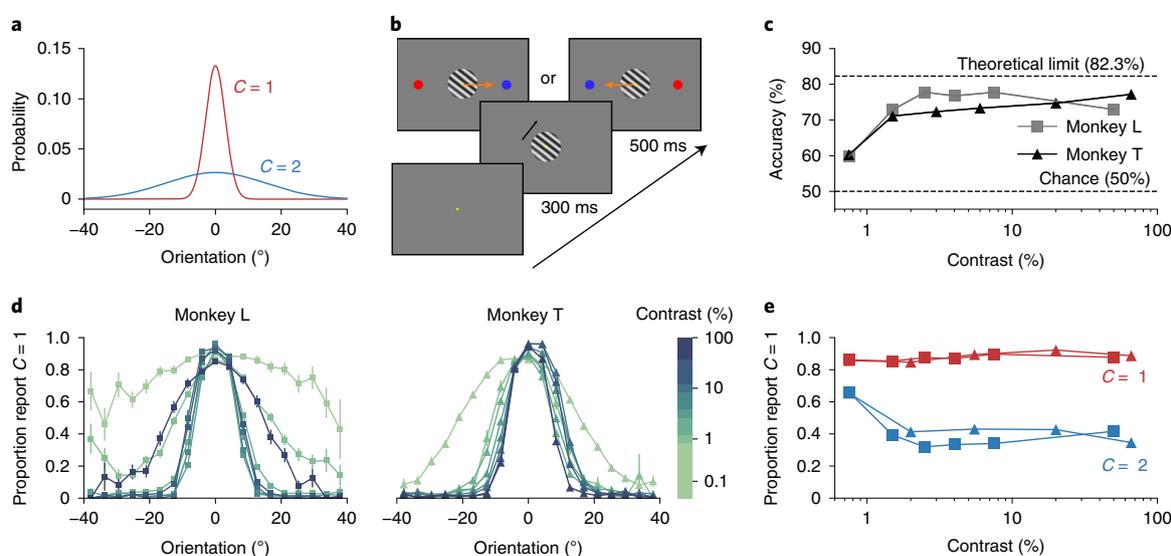


Fig. 2 | Behavioral task. **a**, The stimulus orientation distributions for the two classes. The two distributions shared the same mean ($\mu=0^\circ$) but differed in their standard deviations ($\sigma_1=3^\circ$ and $\sigma_2=15^\circ$). **b**, Time course of a single trial. The subject fixated on 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 ($n=110,695$ and $n=192,631$ total trials for monkeys L and T, respectively). All data points are means, and error bars indicate s.e.m. **e**, Class-conditioned responses. For each subject, the proportions of $C=1$ reports are 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**.

the beginning of the session, and on each trial, one of these values was randomly selected.

In our previous study¹⁵, 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 use trial-by-trial uncertainty about the stimulus orientation. In this 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 closely matched the performance of monkey L (Fig. 2c). Both animals had psychometric curves displaying the expected strong dependence on both contrast and orientation (Fig. 2d,e).

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

Decoding likelihood function from V1. Each monkey was implanted with a chronic multielectrode (Utah) array in the parafoveal primary visual cortex (V1) to record the simultaneous cortical population activity as the subjects performed the orientation classification task (Fig. 3a).

A total of 61 and 71 sessions were analyzed from monkeys L and T for a total of 110,695 and 192,631 trials, respectively (Extended Data Fig. 1). In each recording session, up to 96 channels were recorded. On each trial and for each channel, we computed the total number of spikes that occurred during the 500 ms of

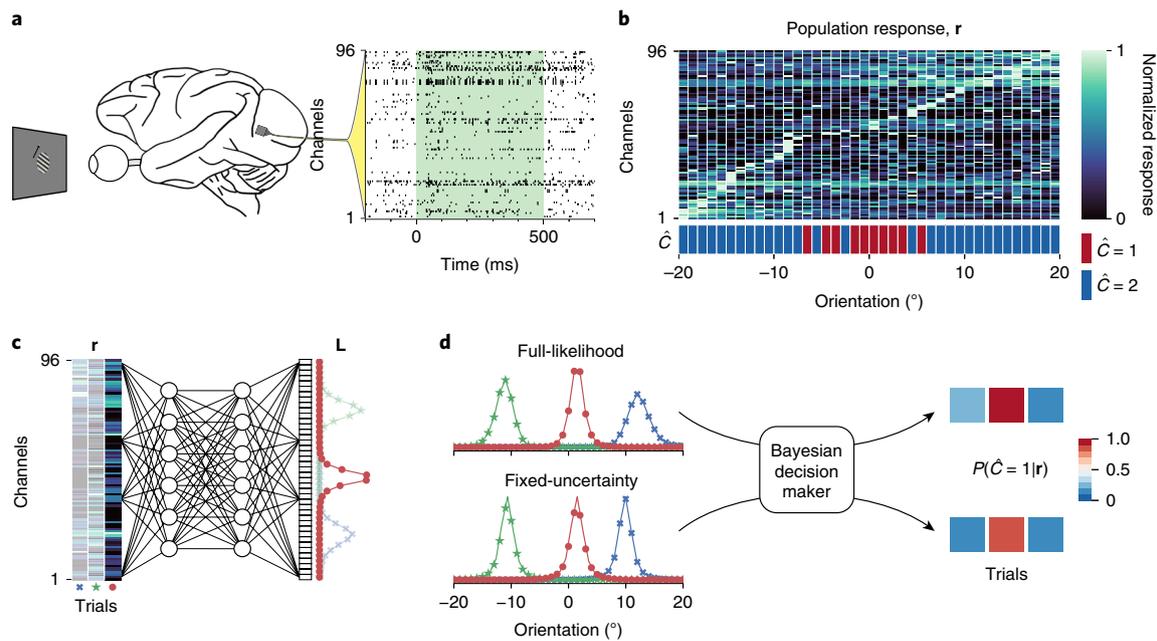


Fig. 3 | Encoding and decoding of the stimulus orientation. **a**, An example of 96-channel spike traces from a single trial (monkey T). The vector of spike counts, \mathbf{r} , was accumulated over the presaccade stimulus presentation period (time 0–500 ms, green shade). **b**, The population response for the selected trials from a single contrast-session (monkey T, 64% contrast). Column: a population response \mathbf{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 \mathbf{r} to the decoded likelihood function, \mathbf{L} . All likelihood functions are area-normalized. **d**, Two decision models M based on alternative likelihood decoders. In the full-likelihood model using the full-likelihood decoder, the likelihood \mathbf{L} was decoded without any constraints on the shape. In the fixed-uncertainty model using 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, and the values of the likelihood functions were read out at the discrete points indicated by the markers. For both models, the resulting likelihood functions were fed into a parameterized Bayesian decision maker to yield the decision prediction $P(\hat{C}=1|\mathbf{r}, M)$.

stimulus presentation preceding the decision-making cue (Fig. 3a), yielding a vector of population responses \mathbf{r} used in the subsequent analyses (Fig. 3b).

Existing computational methods for decoding the trial-by-trial likelihood function from the cortical population activities typically make strong parametric assumptions about the stimulus-conditioned distribution of the population response (that is, the generative model of the population response). For example, population responses to a stimulus can be modeled as an independent Poisson distribution, allowing each recorded unit to be characterized by a simple tuning curve (which may be further parameterized)^{14,18–22}. Although this simplifying assumption makes computing the trial-by-trial likelihood function straightforward, disregarding potential correlations among the units in population responses (that is, noise correlations and internal brain state fluctuations^{23–28}) can lead to biased estimates of the likelihood function and limit the generality of this approach. Even though more generic parametric models, such as Poisson-like distributions, of population distributions have been proposed^{9,10,15,29,30}, 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 allowed 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)¹⁷ to predict the per-trial likelihood function $\mathcal{L}(\theta) \equiv P(\mathbf{r}|\theta)$ over

stimulus orientation θ 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 used in training the DNN.

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

Trial-to-trial uncertainty improves behavioral predictions. To assess whether the uncertainty decoded from population responses in the form of sensory likelihood functions mediates the behavioral outcome (perceptual decision) as we hypothesized, it is critical that we appropriately condition the analysis on the stimulus. To illustrate the importance of conditioning on the stimulus to determine whether the decoded likelihood function mediates perceptual decisions, consider a typical perceptual decision-making task (like ours) (Extended Data 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 (for example, 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 (for example, orientation or contrast) across trials can drive variation both in the decoded likelihood function and in the animal's decision, one may find a dependence of \hat{C} on \mathcal{L} , 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 .

Because 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^{31,32}, 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 that uses 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 earlier (Fig. 3d); and (2) a fixed-uncertainty model that uses 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 (Extended Data Fig. 5). The fixed-uncertainty model captures the alternative hypothesis in which the recorded sensory population encodes only a point estimate of the sensory variable (that is, 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).

We observed that likelihood functions decoded by the full-likelihood decoder exhibited the expected dependencies on the overt drivers of uncertainty such as contrast (Fig. 4a–c): the width of the likelihood function was higher at lower contrast (Fig. 4d), and generally, the likelihood function decoded by the fixed-uncertainty decoder closely approximated the likelihood function decoded by the full-likelihood decoder (Extended Data Fig. 5). We use the term 'decoder' for the DNN that returns estimated likelihood functions and the term 'decision maker' for the mapping from likelihood function to decision. We refer to the combination of a decoder and a decision maker as the 'decision model', M .

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 (Extended Data Fig. 6; refer to Methods for details). The model parameters were fitted by maximizing the total log likelihood over all trials (indexed by i) 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 orientation (for example, the mean of the likelihood function) and differed only 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 used 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 (Extended Data Fig. 7), reaching up to a 90% accuracy rate. 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; 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 models: monkey L, two-tailed paired t test, $t(110694) = 11.06$, $P < 10^{-9}$, $\delta_{\text{total}} = 11.0 \times 10^2$ over a total of 110,695 trials; and monkey T, $t(192630) = 11.03$, $P < 10^{-9}$, $\delta_{\text{total}} = 11.3 \times 10^2$ over a total of 192,631, where δ_{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 250 ms of stimulus presentation. We found a similar improvement for the full-likelihood model over the fixed-uncertainty model in both periods (Extended Data Fig. 8).

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 upper bound on the total log-likelihood differences (δ_{total}) 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 s.d. across five repetitions of simulation for monkeys L and T, respectively) were larger but in the same order of the magnitude as the observed model performance differences (11.0×10^2 and 11.3×10^2 total log-likelihood differences, δ_{total} , across all trials for monkeys 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 (Extended Data Fig. 9).

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 (Extended Data Fig. 9). By design, the fixed-uncertainty model cannot capture this stimulus-dependent change in uncertainty, which could contribute

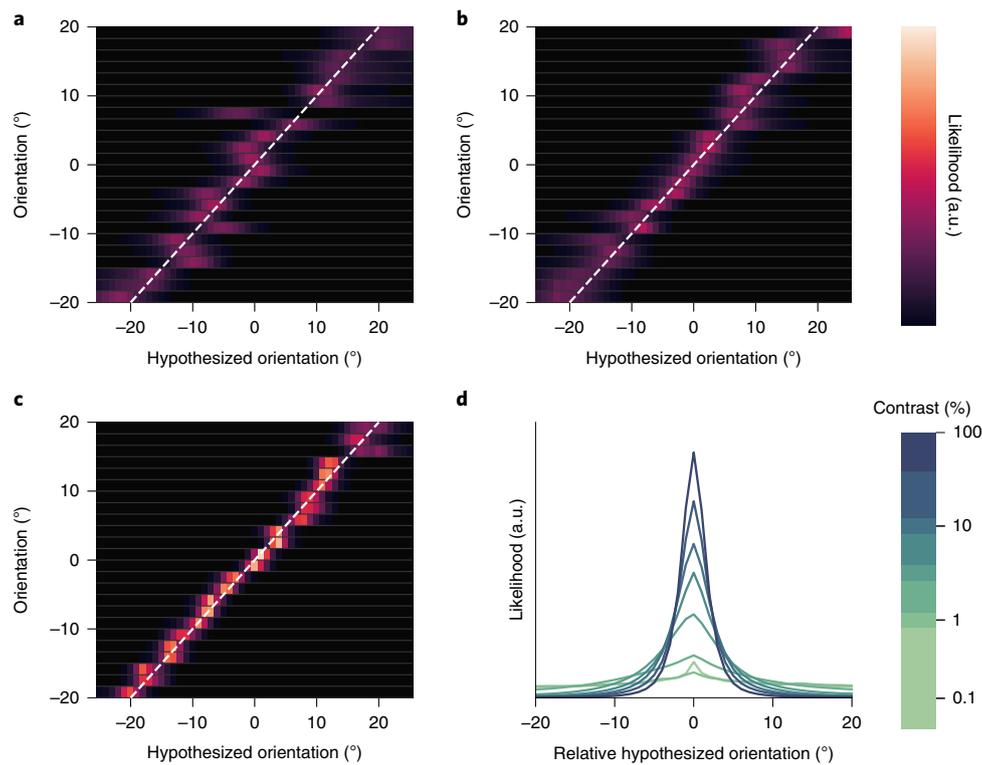


Fig. 4 | Likelihood functions decoded by the trained full-likelihood decoders. a–c, Example decoded likelihood functions from three contrast-sessions, 4% (**a**), 32% (**b**) and 100% (**c**), from monkey T. Each row represents the decoded likelihood function over the hypothesized orientation for a randomly 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°. The centered likelihood functions were then averaged across all trials within the same contrast bin. a.u., arbitrary units.

to it performing worse than the full-likelihood model (Extended Data 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 preserved 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, two-tailed paired t test $t(110694) = -18.44$, $P < 10^{-9}$, $\delta_{\text{total}} = -20.9 \times 10^2$ over a total of 110,695 trials; and monkey T: $t(192630) = -20.15$, $P < 10^{-9}$, $\delta_{\text{total}} = -25.9 \times 10^2$ over a total of 192,631 trials, where δ_{total} is the total log-likelihood difference across all trials). Therefore, there were significant performance differences in the full-likelihood model

between the unshuffled and shuffled data (trial log-likelihood difference: monkey L, two-tailed paired t test $t(110,694) = 33.34$, $P < 10^{-9}$, $\delta_{\text{total}} = 31.9 \times 10^2$; and monkey T, $t(192,630) = 34.52$, $P < 10^{-9}$, $\delta_{\text{total}} = 37.2 \times 10^2$).

To confirm that our effect sizes were appropriate, we again compared these values with those obtained from simulations in which we took the full-likelihood model to be the ground truth (Extended Data Fig. 9). 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 s.d. across five repetitions), which was similar in magnitude to the observed δ_{total} 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 (that is, mean of the likelihood function) also encoded the uncertainty regarding that estimate (that is, shape of the likelihood function), as we hypothesized to be the case, we performed attribution analysis³³ on the trained full-likelihood decoder. Through this analysis, we ask how much of the changes in either the mean of the likelihood μ_L (that is, surrogate for the best point estimate) or the standard deviation of the likelihood function σ_L (that is, surrogate measure of the uncertainty) can be attributed back to each input multiunit, yielding attribution vectors \mathbf{A}_μ and \mathbf{A}_σ , respectively. The question of feature attribution is an active field of research in machine learning, and multiple methods

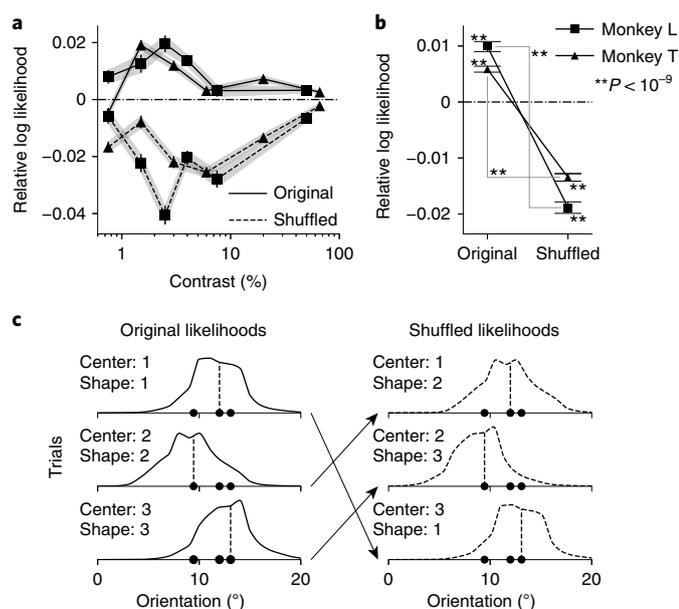


Fig. 5 | Model performance. **a**, Average trial-by-trial performance of the full-likelihood model relative to the fixed-uncertainty model across contrasts, measured as the average trial difference in the log likelihood computed over total trials of $n = 1,110,695$ and $n = 192,631$ for monkeys L and T, respectively. The results for the original (unshuffled) and the shuffled data are shown by solid and dashed lines, respectively. The squares and triangles mark monkeys 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 trial log-likelihood difference between the full-likelihood and fixed-uncertainty models was significant for both monkeys on both the original (two-tailed paired t test: $t(110,694) = 11.06$ with $P < 10^{-9}$ for monkey L; $t(192,630) = 11.03$ with $P < 10^{-9}$ for monkey T) and the shuffled data (two-tailed paired t test: $t(110,694) = -18.44$ with $P < 10^{-9}$ for monkey L; $t(192,630) = -20.15$ with $P < 10^{-9}$ for monkey T). Furthermore, the difference between the full-likelihood model on the original and the shuffled data was significant (two-tailed paired t test: $t(110,694) = 33.34$ with $P < 10^{-9}$ for monkey L; $t(192,630) = 34.52$ with $P < 10^{-9}$ for monkey T). **a, b**, All data points are means, and error bars and shaded areas indicate s.e.m. All P values are Bonferroni-corrected for the three comparisons for each monkey. **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.

of attribution computation exist^{33–35}. Here, we have selected three different methods of computing attribution: saliency maps³⁴, gradient \times input³³ and DeepLIFT (DL)³⁵ (refer to Methods for the details of attribution computation).

We observed that across all three attribution methods, multiunits with high μ_L attribution tended to have high σ_L attribution, and vice versa, giving rise to a high degree of correlation between elements of A_μ and A_σ , denoted as A_μ and A_σ (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 μ_L attribution to have a low σ_L attribution, and vice versa, therefore leading to negative correlation between A_μ and A_σ . However, we observed that across all contrast-sessions from both monkeys, A_μ was strongly positively correlated with A_σ 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). We further observed that a substantial fraction of multiunits participated

in encoding both μ_L (59, 28 and 38 out of 96 multiunits needed to attain >90% of total attribution under saliency maps, gradient \times input and DL, respectively) and σ_L (60, 28 and 38 out of 96 multiunits needed to attain >90% of total attribution under saliency maps, gradient \times input and DL, respectively), suggesting that the information about the encoded likelihood functions is distributed across neurons rather than being encoded by only a small number of neurons in the population (Extended Data Fig. 10).

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 and Kulikowski³⁶ and Britten et al.³⁷ Here, we go beyond this literature by examining the hypothesis that the brain takes into account knowledge of the form of neural variability 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-making model (the Bayesian decision maker) to predict behavior. We found that a decision 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 population-level evidence in support of the theoretical framework of PPC, 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 DNN 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 (for example, contrast, blur, motion coherence) to manipulate uncertainty^{7,8,22,38}. 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 (Extended Data 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 , that is, $\log \mathcal{L}(s) = \log P(\mathbf{r}|s) + b(\mathbf{r})$ as a function of s . We never commit to a particular generative model $P(\mathbf{r}|s)$ as a function of \mathbf{r} , because the DNN has an arbitrary offset as a function of \mathbf{r} (equation (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 (Extended Data Fig. 7). That being said, the DNN can use what we know about generative models in

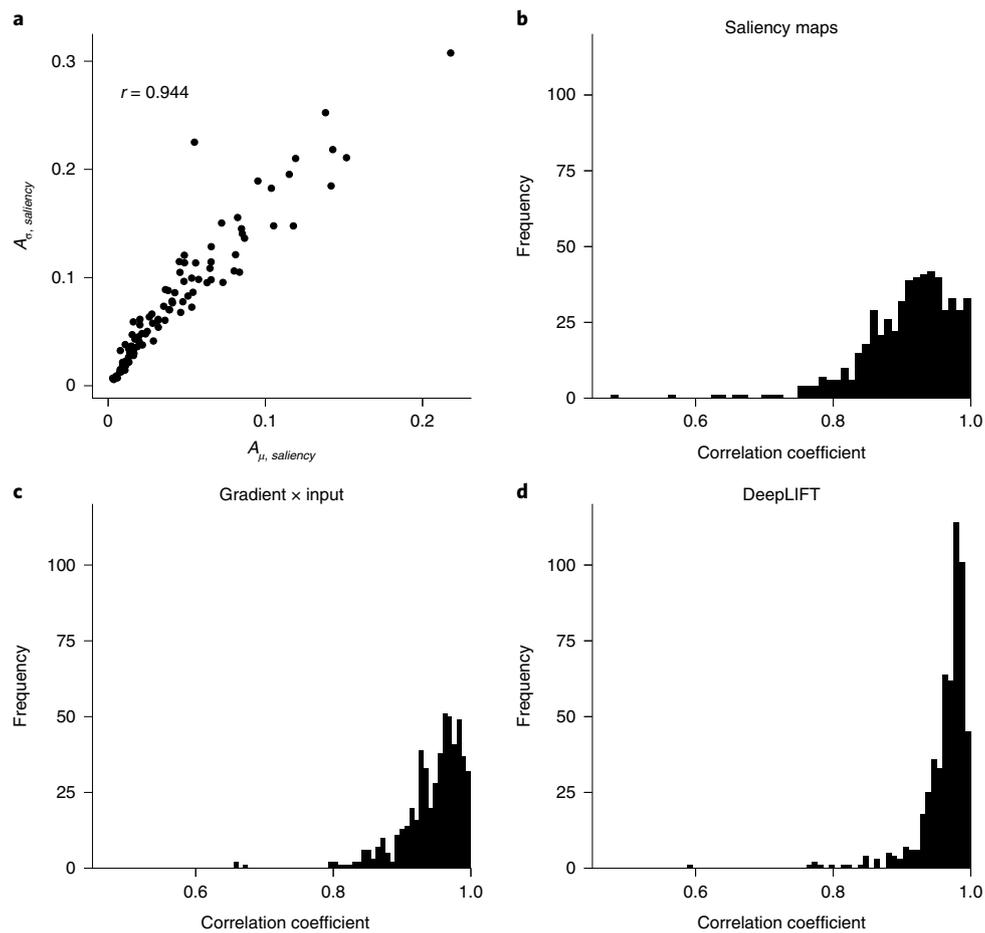


Fig. 6 | Attribution analysis for means and standard deviations of the likelihood functions. **a**, Attribution of 96 input multiunits to the likelihood mean $A_{\mu, \text{saliency}}$ versus s.d. The $A_{\sigma, \text{saliency}}$ value was computed based on saliency maps for an example contrast-session (monkey T, 32% contrast). The value of Pearson's correlation coefficient is given as r in the plot. **b–d**, Distribution of Pearson's correlation coefficients between A_{μ} and A_{σ} for multiunits across all contrast-sessions ($n = 197$ and $n = 349$ total contrast-sessions for monkeys L and T, respectively), which was computed based on different attribution methods.

visual cortex (for example, tuning curves, contrast gain) and also allows for rich correlation among units in the population. Third, interpreting what the DNN does requires more sophisticated deep learning-based techniques. A widely studied question in population coding is how much individual neurons contribute to the population-level information as measured through choice probabilities³⁷. The features of our task that make it sensitive to the use of full likelihood functions also make it brittle for estimating choice probabilities, simply because there are few trials at the psychophysical 50/50 point. Our attribution analysis offers a partial remedy, because it shows that a large number of multiunits drive both the mean orientation and the likelihood width from the decoder (Fig. 6a and Extended Data Fig. 10). However, because the decoder was trained only on visual stimuli and not on behavior, directly linking the attribution analysis to choice probability is challenging. Fourth, we would like to stress that we do not believe that the DNN that we used to decode the likelihood is literally implemented in the brain. What kind of transformation the brain performs to use and compute with this information remains an important question and avenue for future research.

Although 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 downstream

computation that relied heavily on the assumption of Poisson-like variability^{9,10,15,39–41}. However, more recent work has demonstrated that training generic shallow networks accomplishes the same goal without the need for task-specific manual construction⁴². 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^{28,43,44}; in this work, we are agnostic to this issue. Although 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^{45,46} to process models whose internal building blocks, likelihood functions and probability distributions, are more concretely instantiated in neuronal activity^{6,47,48}.

Online content

Any methods, additional references, Nature Research reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author

contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41593-019-0554-5>.

Received: 7 October 2018; Accepted: 6 November 2019;
Published online: 23 December 2019

References

- Laplace, P.-S. *Theorie Analytique des Probabilites* (Ve Courcier, Paris, 1812).
- von Helmholtz, H. Versuch einer erweiterten Anwendung des Fechnerschen Gesetzes im farbensystem. *Z. Psychol. Physiol. Sinnesorg* **2**, 1–30 (1891).
- Knill, D. C. & Richards, W. (eds) *Perception as Bayesian Inference* (Cambridge University Press, 1996).
- Kersten, D., Mamassian, P. & Yuille, A. Object perception as Bayesian inference. *Annu. Rev. Psychol.* **55**, 271–304 (2004).
- Knill, D. C. & Pouget, A. The Bayesian brain: the role of uncertainty in neural coding and computation. *Trends Neurosci.* **27**, 712–719 (2004).
- Ma, W. J. & Jazayeri, M. Neural coding of uncertainty and probability. *Annu. Rev. Neurosci.* **37**, 205–220 (2014).
- Alais, D. & Burr, D. The ventriloquist effect results from near-optimal bimodal integration. *Curr. Biol.* **14**, 257–262 (2004).
- Ernst, M. O. & Banks, M. S. Humans integrate visual and haptic information in a statistically optimal fashion. *Nature* **415**, 429–433 (2002).
- Ma, W. J., Beck, J. M., Latham, P. E. & Pouget, A. Bayesian inference with probabilistic population codes. *Nat. Neurosci.* **9**, 1432–1438 (2006).
- Beck, J. M. et al. Probabilistic population codes for bayesian decision making. *Neuron* **60**, 1142–1152 (2008).
- Pouget, A., Dayan, P. & Zemel, R. Information processing with population codes. *Nat. Rev. Neurosci.* **1**, 125–132 (2000).
- Pouget, A., Dayan, P. & Zemel, R. S. Inference and computation with population codes. *Annu. Rev. Neurosci.* **26**, 381–410 (2003).
- Ma, W. J., Beck, J. M. & Pouget, A. Spiking networks for Bayesian inference and choice. *Curr. Opin. Neurobiol.* **18**, 217–222 (2008).
- Graf, A. B. A., Kohn, A., Jazayeri, M. & Movshon, J. A. Decoding the activity of neuronal populations in macaque primary visual cortex. *Nat. Neurosci.* **14**, 239–245 (2011).
- Qamar, A. T. et al. Trial-to-trial, uncertainty-based adjustment of decision boundaries in visual categorization. *Proc. Natl. Acad. Sci. USA* **110**, 20332–20337 (2013).
- LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436–444 (2015).
- Goodfellow, A., Bengio, I. & Courville, Y. *Deep Learning* (MIT Press, 2016).
- Seung, H. S. & Sompolinsky, H. Simple models for reading neuronal population codes. *Proc. Natl. Acad. Sci. USA* **90**, 10749–10753 (1993).
- Sanger, T. D. Probability density estimation for the interpretation of neural population codes. *J. Neurophysiol.* **76**, 2790–2793 (1996).
- Zemel, R. S., Dayan, P. & Pouget, A. Probabilistic interpretation of population codes. *Neural Comput.* **10**, 403–430 (1998).
- Jazayeri, M. & Movshon, J. A. Optimal representation of sensory information by neural populations. *Nat. Neurosci.* **9**, 690–696 (2006).
- Fetsch, C. R., Pouget, A., DeAngelis, G. C. & Angelaki, D. E. Neural correlates of reliability-based cue weighting during multisensory integration. *Nat. Neurosci.* **15**, 146–154 (2012).
- Averbeck, B. B. & Lee, D. Effects of noise correlations on information encoding and decoding. *J. Neurophysiol.* **95**, 3633–3644 (2006).
- Ecker, A. S. et al. Decorrelated neuronal firing in cortical microcircuits. *Science* **327**, 584–587 (2010).
- Ecker, A. S., Berens, P., Tolias, A. S. & Bethge, M. The effect of noise correlations in populations of diversely tuned neurons. *J. Neurosci.* **31**, 14272–14283 (2011).
- Ecker, A. S. et al. State dependence of noise correlations in macaque primary visual cortex. *Neuron* **82**, 235–248 (2014).
- van Bergen, R. S. & Jehee, J. F. M. Modeling correlated noise is necessary to decode uncertainty. *Neuroimage* **180**, 78–87 (2018).
- Denfield, G. H., Ecker, A. S., Shinn, T. J., Bethge, M. & Tolias, A. S. Attentional fluctuations induce shared variability in macaque primary visual cortex. *Nat. Commun.* **9**, 2654 (2018).
- Ma, W. J. Signal detection theory, uncertainty, and poisson-like population codes. *Vis. Res.* **50**, 2308–2319 (2010).
- Van Bergen, R. S., Ma, W. J., Pratte, M. S. & Jehee, J. F. M. Sensory uncertainty decoded from visual cortex predicts behavior. *Nat. Neurosci.* **18**, 1728–1730 (2015).
- Tolhurst, D. J., Movshon, J. A. & Dean, A. F. The statistical reliability of signals in single neurons in cat and monkey visual cortex. *Vis. Res.* **23**, 775–785 (1983).
- Shadlen, M. N. & Newsome, W. T. The variable discharge of cortical neurons: implications for connectivity, computation, and information coding. *J. Neurosci.* **18**, 3870–3896 (1998).
- Ancona, M., Ceolini, E., Öztireli, C. & Gross, M. A unified view of gradient-based attribution methods for deep neural networks. In *NIPS 2017 Workshop on Interpreting, Explaining and Visualizing Deep Learning* <http://www.interpretable-ml.org/nips2017workshop/papers/02.pdf> (2017).
- Simonyan, K., Vedaldi, A. & Zisserman, A. Deep inside convolutional networks: visualising image classification models and saliency maps. Preprint at *arXiv* <https://arxiv.org/abs/1312.6034> (2013).
- Shrikumar, A., Greenside, P. & Kundaje, A. Learning important features through propagating activation differences. in *Proceedings of the 34th International Conference on Machine Learning, Proceedings of Machine Learning Research* Vol. 70 (eds Precup, D. & Teh, Y. W.) 3145–3153 (2017).
- Campbell, F. W. & Kulikowski, J. J. The visual evoked potential as a function of contrast of a grating pattern. *J. Physiol.* **222**, 345–356 (1972).
- 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. *Vis. Neurosci.* **13**, 87–100 (1996).
- Angelaki, D. E., Humphreys, G. & DeAngelis, G. C. Multisensory integration. *J. Theor. Humanit.* **19**, 452–458 (2009).
- Ma, W. J., Navalpakkam, V., Beck, J. M., van den Berg, R. & Pouget, A. Behavior and neural basis of near-optimal visual search. *Nat. Neurosci.* **14**, 783–790 (2011).
- Beck, J. M., Latham, P. E. & Pouget, A. Marginalization in neural circuits with divisive normalization. *J. Neurosci.* **31**, 15310–15319 (2011).
- Ma, W. J. & Rahmati, M. Towards a neural implementation of causal inference in cue combination. *Multisens. Res.* **26**, 159–176 (2013).
- Orhan, A. E. & Ma, W. J. Efficient probabilistic inference in generic neural networks trained with non-probabilistic feedback. *Nat. Commun.* **8**, 138 (2017).
- Cumming, B. G. & Nienborg, H. Feedforward and feedback sources of choice probability in neural population responses. *Curr. Opin. Neurobiol.* **37**, 126–132 (2016).
- Bondy, A. G., Haefner, R. M. & Cumming, B. G. Feedback determines the structure of correlated variability in primary visual cortex. *Nat. Neurosci.* **21**, 598–606 (2018).
- Geisler, W. S. Contributions of ideal observer theory to vision research. *Vis. Res.* **51**, 771–781 (2011).
- Körding, K. Decision theory: what 'should' the nervous system do? *Science* **318**, 606–610 (2017).
- Maloney, L. T. & Mamassian, P. Bayesian decision theory as a model of human visual perception: testing Bayesian transfer. *Vis. Neurosci.* **26**, 147–155 (2009).
- Ma, W. J. Organizing probabilistic models of perception. *Trends Cogn. Sci.* **16**, 511–518 (2012).

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

© The Author(s), under exclusive licence to Springer Nature America, Inc. 2019

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 weighing 9.5 and 15.1 kg, respectively. All experimental procedures complied with guidelines of the National Institutes of Health 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 12-h light–dark cycle, together with about 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.

No statistical methods were used to predetermine sample sizes, but our sample sizes are similar to those reported in previous publications^{26,28}. Data collection and analysis were not performed blind to the conditions of the experiments. In performing the analysis, no animal was excluded from the study. Some data points (that is, trials) were excluded from the analysis based on criteria described in detail later in this article (see Dataset and inclusion criteria). Additional details may be found in Life Sciences Reporting Summary.

Stimulus presentation. Each visual stimulus was a single drifting oriented sinusoidal grating (spatial frequency: 2.79 cycles per degree visual angle, drifting speed: 3.89 cycles s⁻¹) 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 monkeys L and 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 cathode-ray tube monitor (at a distance of 100 cm; resolution: 1,600 × 1,200 pixels; refresh rate: 100 Hz) using Psychophysics Toolbox⁴⁹. The monitor was gamma corrected to have a linear luminance response profile. Video cameras (DALSA genie HM640; frame rate 200 Hz) 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 θ , drawn from one of two classes, each defined by a Gaussian probability distribution. Both distributions have a mean of 0° (grating drifting horizontally rightward, positive orientation corresponding to counterclockwise 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 basis¹⁵.

Each trial proceeded as follows: a trial was initiated by a beeping sound and the appearance of a fixation target (0.15° visual angle) in the center of the screen. The monkey fixated on the fixation target for 300 ms within 0.5°–1° visual angle. The stimulus then appeared at the center of the screen. After 500 ms, two colored targets (red and green) appeared to the left and the right of the grating stimulus (horizontal offset of 4.29° from the center with the target diameter of 0.71° visual angle), at which point the monkey saccaded to one of the targets to indicate its choice of class. For monkey L, the grating stimulus was removed from the screen when the saccade target appeared, whereas 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 was 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^{28,50,51}. In brief, 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⁻¹ at 48 h, 24 h and on the day of the procedure) and atropine (0.05 mg kg⁻¹

before sedation), animals were sedated with a mixture of ketamine (10 mg kg⁻¹) and xylazine (0.5 mg kg⁻¹). 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. Analgesics were given for 7 d after surgery to alleviate pain.

Electrophysiological recording and data processing. The neural signals were preamplified 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, USA) 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^{26,28,50}. In brief, a spike was detected when the signal on a given electrode crossed a threshold of five times the s.d. 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 s.d.⁵², 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/ephys-preprocessing>. In this study, the term ‘multiunit’ refers to the set of all spikes detected from a single channel (that is, 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 pretarget stimulus presentation, r_i for the i 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 monkeys 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 (for example, 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 $r_i = \sum_i r_i$ 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. σ_{adj} is the s.d. 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 th unit spike count r_i for the trial should fall within a range defined as: $|r_i - \text{MED}_i| \leq 1.5 \times \text{IQR}_i$, where MED_i and IQR_i are the median and IQRs of the i th unit spike counts distribution throughout the session, respectively.

We included only trials that satisfied both of the criteria in our analysis. Empirically, we found the earlier 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 monkeys L and T, thus retaining 98.77% and 99.48% of the total trials, respectively. Although this selection criteria allowed us to remove any apparent anomaly in the data, we found that the main findings described in this article were not sensitive to the precise definition of the inclusion criteria.

When analyzing the population responses to the first half (0–250 ms) or the second half (250–500 ms) of the stimulus presentation as found in Extended Data Fig. 8, we reapplied the earlier trial selection criteria to the respective time segment, obtaining 110,816 and 192,962 trials for monkeys L and T for the first half (0–250 ms), and 110,887 and 192,980 trials for monkeys L and T for the second half (0–500 ms), again retaining >98% of the total trials in all conditions.

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° of visual angle presented on a uniform gray background, with randomly varying location and color (black or white) every 30 ms for 1 s. 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 θ , 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 θ . This may be computed through the knowledge of the generative relation leading from θ 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 distribution¹⁹. Under this approach, the expected firing rates (that is, tuning curves) of the i th neuron $E[r_i|\theta] = f_i(\theta)$ must be approximated as well, for example, by fitting a parametric function (for example, von Mises tuning curves⁵³) or employing kernel regression¹⁹. Although these approaches have proved 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 conditional $\mathcal{L}(\theta)$ without having to make strong assumptions about the underlying conditional probability distribution over \mathbf{r} .

To this end, we used a DNN¹⁶ to directly approximate the likelihood function over the stimulus orientation, θ , 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 θ into n distinct values, θ_1 to θ_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, $f: \mathbb{R}^m \mapsto \mathbb{R}^n$:

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

for some scalar function $b \in \mathbb{R}$. As the experimenter, we know the distribution of the stimulus orientation, $\mathbf{p}_\theta \in \mathbb{R}^n$, where $\mathbf{p}_{\theta_i} = P(\theta_i)$. We combine $f(\mathbf{r})$ and \mathbf{p}_θ to compute the log posterior over stimulus orientation θ up to some scalar value $b'(\mathbf{r})$,

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

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

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

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

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

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

The goal then is to train the DNN $f(\mathbf{r})$ such that the overall function $\mathbf{q}(\mathbf{r})$ matches the 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, in turn, allows the network output $f(\mathbf{r})$ to approach the log of the likelihood function \mathbf{L} , up to a constant $b(\mathbf{r})$. For one-out-of- n classification problems, minimizing the cross-entropy between $\mathbf{q}(\mathbf{r})$ and the stimulus orientation θ for a given \mathbf{r} lets the overall function $\mathbf{q}(\mathbf{r})$ approach the true posterior $\mathbf{p}(\mathbf{r})$, as desired^{54,55}. To show this, let us start by minimizing the difference between the model estimated posterior $\mathbf{q}(\mathbf{r})$ and the true posterior $\mathbf{p}(\mathbf{r})$ over the distribution of \mathbf{r} . We do this by minimizing the loss, L , defined as the expected value of the Kullback–Leibler (KL) divergence⁵⁶ between the two posteriors:

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

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

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

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

where $P(\theta = \theta_i|\mathbf{r}) \equiv \mathbf{p}_i(\mathbf{r})$, $q(\theta = \theta_i|\mathbf{r}, W) \equiv \mathbf{q}_i(\mathbf{r}, W)$, W is a collection of all trainable parameters in the network and $H(\theta|\mathbf{r})$ is the conditional entropy of θ conditioned on \mathbf{r} , which is an unknown but a fixed quantity with respect to W and the data distribution. Here, we used the notation $\mathbf{q}(\mathbf{r}, W)$ to highlight the dependence of the network estimated posterior $\mathbf{q}(\mathbf{r})$ on the network parameters W . We can redefine the loss, L^* , leaving only the terms that depend on the trainable parameters W , and then apply a Monte Carlo method⁵⁷ to approximate the loss from N samples:

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

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

where $(\theta^{(i)}, \mathbf{r}^{(i)})$ are samples drawn from a training set for the network.

Equation (11) is precisely the definition of the cross-entropy as we set out to show it.

Therefore, by optimizing the overall function $\mathbf{q}(\mathbf{r})$ to match the posterior distribution through the use of cross-entropy loss, the network output $f(\mathbf{r})$ can approximate the log of the likelihood 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 θ to \mathbf{r} , $P(\mathbf{r}|\theta)$.

As an example, consider a neural population with responses that follow a Poisson-like distribution (that is, a version of the exponential distribution with linear sufficient statistics^{9,10}). 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 $P(\mathbf{r}|\theta) = \phi(\mathbf{r}) e^{h^T(\theta)\mathbf{r}}$. The log-likelihood function is then $\log \mathbf{L} = \log \phi(\mathbf{r}) + \mathbf{H}^T \mathbf{r}$, where \mathbf{H} is a matrix whose i th column is $h(\theta_i)$. If we let $f(\mathbf{r}) = \mathbf{H}^T \mathbf{r}$, 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 $h(\theta)$ of the Poisson-like distribution.

In this work, we modeled the mapping $f(\mathbf{r})$ as a DNN with two hidden layers¹⁷, consisting of two repeating blocks of a fully connected layer of size N_h , followed by a rectified linear unit (ReLU)¹⁶ and a dropout layer⁵⁸ with dropout rate d_p , 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 $h = [-0.25, 0.5, -0.25]$. Therefore, the training loss included the term:

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

for $\mathbf{u} = (\log \mathbf{L}) * h$, where $*$ denotes convolution operation, \mathbf{u}_i is the i th element of the filtered log-likelihood function \mathbf{u} and γ 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 θ 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 monkeys L and T data, respectively. The network was trained on the training set, starting with initial learning rate of λ_0 , and its performance on the validation set was monitored to perform early stopping⁵⁹ 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 θ 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 four times, therefore training the network under the four sequentially decreasing learning rate schedules 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 λ_0 , the weight on the likelihood function smoothness regularizer γ and the dropout rate d_p , 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.

Decoding likelihood functions from known response distributions. To assess the effectiveness of the DNN-based likelihood decoding method described earlier, we simulated neural population responses with known noise distributions, trained DNN decoders on the simulated population responses and compared the decoded likelihood functions with 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 with those decoded by assuming independent Poisson distribution on the population responses, as done in previous work^{14,18,19,21,22}.

We simulated the activities of a population of 96 multiunits \mathbf{r}_{sim} responding to the stimulus orientation θ drawn from 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 th unit to θ —that is, the tuning function $f_i(\theta)$ —with a Gaussian function:

$$\mathcal{L}(\theta) = \prod_i P(r_i|\theta) = \prod_i \frac{f_i(\theta)^{r_i} e^{-f_i(\theta)}}{r_i!} \quad (14)$$

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

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

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

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

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

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

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

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

for

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

where $\mathbf{f}^{1/2}(\theta) \in \mathbb{R}^{96}$, such that $f_i^{1/2}(\theta) = \sqrt{f_i(\theta)}$ and $\mathbf{C} \in \mathbb{R}^{96 \times 96}$ 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 \mathbf{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 1,200 trials. Among these, 200 trials were set aside as the test set. We trained the DNN-based likelihood decoder on the remaining 1,000 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}_{sim} under the assumption of independent Poisson variability, regardless of the 'true' distribution. Each unit's responses over the 1,000 trials were fitted separately with a Gaussian tuning curve (equation (14)). The parameters of the tuning curve A , μ_i and $\sigma_{\text{sim},i}$ were obtained by maximizing the likelihood of observing the i th unit's responses ($\theta, r_{\text{sim},i}$) for the given Gaussian tuning curve assuming independent Poisson noise, using `minimize` function from the 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 equation (16) for the independent Poisson population or equation (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}_{DNN} by computing their KL divergence⁵⁶ from the ground-truth likelihood function $\mathcal{L}_{\text{gt}}(\theta)$, giving rise to D_{Pois} and D_{DNN} , respectively. All continuous likelihood functions (\mathcal{L}_{gt} and $\mathcal{L}_{\text{Pois}}$) were sampled at orientation θ , where $\theta \in \mathbb{Z}$ and $\theta \in [-45^\circ, 45^\circ]$, giving rise to discretized likelihood functions \mathbf{L}_{gt} and \mathbf{L}_{Pois} , matching the dimensionality of the discretized likelihood function \mathbf{L}_{DNN} computed by the DNN. We then computed the KL divergence as:

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

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 (Extended Data Fig. 3). When the simulated population distribution was independent Poisson, then $D_{\text{Pois}} < D_{\text{DNN}}$ for

all test set trials, indicating that \mathbf{L}_{Pois} better approximated \mathbf{L}_{gt} overall than \mathbf{L}_{DNN} . However, \mathbf{L}_{DNN} still closely approximated \mathbf{L}_{gt} .

When the simulated population distribution was mean scaled correlated Gaussian, \mathbf{L}_{DNN} better approximated \mathbf{L}_{gt} than \mathbf{L}_{Pois} on the majority of the trials. Furthermore, \mathbf{L}_{Pois} provided qualitatively worse fit to the \mathbf{L}_{gt} for the simulated correlated Gaussian distribution compared with the fit of \mathbf{L}_{DNN} to \mathbf{L}_{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 the 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 a reasonable estimate of the likelihood function across a 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, \mathbf{L}_0 , and a DNN that reads out a single scalar value, Δ , corresponding to the shift that is applied to \mathbf{L}_0 (Extended Data Fig. 5) to yield the final likelihood function \mathbf{L} . The DNN consisted of two repeating blocks of a fully connected layer followed by ReLU and a dropout layer, and a final fully connected readout layer with no output nonlinearity, much like the DNN used for the full-likelihood decoder. The log \mathbf{L}_0 was shifted by Δ , using linear interpolation-based grid sampling⁶⁰ 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 Δ (and therefore to the DNN parameters) and to the likelihood function shape \mathbf{L}_0 .

The output shifted log-likelihood function was then trained in an identical manner to the full-likelihood decoder described earlier, using 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 $P(\mathbf{r}|\theta)$ (Extended Data Fig. 2).

As noted earlier, decoding likelihood function under the Poisson-like distribution is a special case of the full-likelihood decoder but using entirely linear DNN (that is, no nonlinearity used 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 dropout layer⁵⁸, 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 the likelihood function under the independent Poisson distribution assumption, we first fitted tuning curves $f_i(\theta)$ for each multiunit's responses to stimulus orientations θ within a single contrast-session. Tuning curves were computed using Gaussian process regression⁶¹ with squared exponential covariance function $\text{cov}(f(\theta_1), f(\theta_2)) = \exp\left(-\frac{1}{2\sigma_L}(\theta_1 - \theta_2)^2\right)$ and a fixed observational noise σ_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 s.d. 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 μ_L and σ_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 the precise choice of the point estimate does not affect the main findings.

Attribution analysis. To assess whether the same members of the population simultaneously encode the best point estimate (that is, in the form of the mean of the likelihood function μ_L) and uncertainty (that is, in the form of the width of the likelihood function σ_L), we computed the attribution of each multiunit input of the trained full-likelihood decoder to the mean of the likelihood μ_L and the s.d. of the likelihood function σ_L , giving rise to the attribution $\mathbf{A}_\mu, \mathbf{A}_\sigma \in \mathbb{R}^m$, respectively, where m is the number of multiunits in the input to the network. Among numerous methods of computing attribution^{33–35,62}, we have selected three popular gradient-based attribution methods³³, saliency maps³⁴, gradient \times input⁶² and DL³⁵, and compared their results.

Given a collection of input population responses and computed likelihood functions $\{\mathbf{r}^{(k)}, \mathbf{L}^{(k)}\}$, where the superscript denotes the k th trial in the contrast-session, we compute the mean and the s.d. of the likelihood function according to equations (23) and (24), respectively, giving rise to $\mu_L^{(k)}$ and $\sigma_L^{(k)}$. 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 $\mathbf{a} \in \mathbb{R}^m$. The sign of the attribution indicates whether increasing the unit tends to increase or decrease the target feature. Because we are interested more in how much each unit contributes 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:

$$\mathbf{A}_{S, \text{method}} = \frac{1}{N} \sum_k |\mathbf{a}_{S, \text{method}}^{(k)}| \quad (25)$$

where N is the total number of trials in the contrast-session.

For the saliency maps-based method³⁴, 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 the 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} :

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

where \odot denotes the Hadamard (element-wise) product.

Finally, we computed DL attribution by using modified gradient computation for ReLUs 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 Ancona et al.³³ and Shrikumar et al.³⁵ 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 yields the modified partial derivative $\frac{\partial^m S}{\partial \mathbf{r}}$, which is finally used to compute the DL attribution as:

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

For each contrast-session and each attribution method, we computed the attribution of the input units to both μ_L and σ_L , yielding vectors \mathbf{A}_μ and \mathbf{A}_σ , and we computed Pearson's correlation coefficient between their elements, A_μ and A_σ (Fig. 6). Furthermore, for each contrast-session and each attribution method, the input

units were ordered from the largest to the smallest attribution, and the cumulative attribution over the ordered units was computed (Extended Data Fig. 10).

Among 546 total contrast-sessions collected from monkeys L and T, the full-likelihood decoder trained on one of the contrast-sessions (monkey T, 0.05% contrast with 162 trials) failed to show dependency on the input population responses \mathbf{r} , and therefore attribution could not be properly computed. Given this, attribution analyses were performed on the remaining 545 contrast-sessions.

Decision 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 uses 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 $\mu=0^\circ$, $\sigma_1=3^\circ$ and $\sigma_2=15^\circ$, where $\mathcal{N}(\theta; \mu, \sigma^2)$ denotes a Gaussian distribution over θ with mean μ and variance σ^2 . The posterior ratio ρ 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 MAP decision-making:

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

where \hat{C} is the subject's decision. However, according to the posterior probability matching strategy^{63,64}, the decision of subjects on certain tasks is 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, whereas $\alpha=\infty$ corresponds to the MAP strategy⁶⁴. We fitted the value of α for each contrast-session during the model fitting to capture any variation of the strategy. Furthermore, we incorporated a lapse rate λ , 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 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: μ , $P(C=1)$, α and λ . Fitting μ (the center of stimulus orientation distributions) and $P(C=1)$ (prior over class) allowed us to capture the bias in the stimulus (orientation) distribution that the subject may have acquired erroneously during the training, and fitting α and λ allowed for the model to match the decision-making strategy employed by the subject.

Using the likelihood function $\mathcal{L}(\theta)$ decoded from the V1 population response via the full-likelihood decoder network in equation (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, using 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 θ by shifting the learned fixed-shape likelihood function over orientation. For each contrast-session, a 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 using 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, respectively.

Model fitting and model comparison. We used tenfold 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 ten equally sized subsets, $B_1, B_2, \dots, B_p, \dots, B_{10}$, where B_i is the i th subset. We then held out a single subset B_i as the test set and trained the decision model on the remaining nine 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 were then reported. We repeated this ten 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:

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

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

where Θ is the collection of the parameters for the decision model M and \hat{C}_i is the subject's decision on the i th trial in the training set. The term $P(\hat{C}_i | M, \Theta)$ is given by equation (38) with $\mathcal{L}(\theta)$ from the full-likelihood decoder for the full-likelihood model or from the fixed-uncertainty decoder for the fixed-uncertainty model.

The optimizations were performed using three algorithms: `fmincon` and `ga` from MATLAB's optimization toolbox and Bayesian Adaptive Direct Search⁶⁵. 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 (that is, 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 (that is, 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 (that is, bin center $\in \mathbb{Z}$) with the bin width of 1° . 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, we shifted the newly assigned likelihood function such that the mean of the normalized likelihood function, μ_i (equation (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 (that is, μ_i) (Fig. 5c). To ensure that both models can take 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 equation (38), using 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 five times, thereby producing five sets of simulated data. For each set of simulated data, we trained the two decision models (full-likelihood model and fixed-uncertainty model) on each contrast-session with tenfold cross-validation and reported the aggregated test set scores as was done for the original data.

Statistics. All statistical tests used, including statistic values, sample sizes and P values, are provided in the figure captions. Where t test was used, the underlying data distribution was assumed to be normal, but this was not formally tested. Exact P values less than 10^{-9} were reported as $P < 10^{-9}$. When appropriate, P values were corrected for multiple comparisons using Bonferroni correction, and the corrected P value was reported.

Reporting Summary. Further information on research design is available in the Nature Research Reporting Summary linked to this article.

Data availability

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

Code availability

Codes used for modeling and training the DNNs, as well as for figure generation, can be viewed and downloaded from https://github.com/eywalker/v1_likelihood. All other codes used for analysis, including data selection and decision model fitting, can be found at https://github.com/eywalker/v1_project. Finally, codes used for electrophysiology data processing can be found in the Tolias lab GitHub organization website (<https://github.com/atlab>).

References

- Brainard, D. H. The psychophysics toolbox. *Spat. Vis.* **10**, 433–436 (1997).
- Tolias, A. S. et al. Recording chronically from the same neurons in awake, behaving primates. *J. Neurophysiol.* **98**, 3780–3790 (2007).
- Subramanian, M., Ecker, A. S., Berens, P. & Tolias, A. S. Macaque monkeys perceive the flash lag illusion. *PLoS ONE* **8**, e58788 (2013).
- Quiroga, R. Q., Nadasdy, Z. & Ben-Shaul, Y. Unsupervised spike detection and sorting with wavelets and superparamagnetic clustering. *Neural Comput.* **16**, 1661–1687 (2004).
- Kohn, A. & Movshon, J. A. Adaptation changes the direction tuning of macaque MT neurons. *Nat. Neurosci.* **7**, 764–772 (2004).
- Richard, M. D. & Lippmann, R. P. Neural network classifiers estimate bayesian a posteriori probabilities. *Neural Comput.* **3**, 461–483 (1991).
- Kline, D. M. & Berardi, V. L. Revisiting squared-error and cross-entropy functions for training neural network classifiers. *Neural Comput. Appl.* **14**, 310–318 (2005).
- Kullback, S. & Leibler, R. A. On information and sufficiency. *Ann. Math. Stat.* **22**, 79–86 (1951).
- MacKay, D. J. C. *Information Theory, Inference, and Learning Algorithms* Vol. 22 (Cambridge University Press, 2003).
- Srivastava, N., Hinton, G., Krizhevsky, A. & Salakhutdinov, R. Dropout: a simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.* **15**, 1929–1958 (2014).
- Prechelt, L. in *Neural Networks: Tricks of the Trade* (eds Grégoire, M., Orr, G. B. & Müller, K.-R.) 53–68 (Springer-Verlag, 1998).
- Jaderberg, M., Simonyan, K., Zisserman, A. & Kavukcuoglu, K. Spatial transformer networks. *Adv. Neural Inf. Process. Syst.* **28**, 2017–2025 (2015).
- Rasmussen, C. E. & Williams, C. K. I. *Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning)* (The MIT Press, 2005).
- Shrikumar, A., Greenside, P., Shcherbina, A. & Kundaje, A. Not just a black box: learning important features through propagating activation differences. Preprint at *arXiv* <https://arxiv.org/abs/1605.01713> (2016).
- Mamassian, P. & Landy, M. S. Observer biases in the 3D interpretation of line drawings. *Vis. Res.* **38**, 2817–2832 (1998).
- Acerbi, L., Vijayakumar, S. & Wolpert, D. M. On the origins of suboptimality in human probabilistic inference. *PLoS Comput. Biol.* **10**, e1003661 (2014).

65. Acerbi, L. & Ma, W. J. Practical Bayesian optimization for model fitting with Bayesian adaptive direct search. *Adv. Neural Inf. Process. Syst.* **30**, 1836–1846 (2017).

Acknowledgements

The research was supported by a National Science Foundation Grant (no. IIS-1132009 to W.J.M. and A.S.T.), a DP1 EY023176 Pioneer Grant (to A.S.T.) and grants from the US Department of Health & Human Services, National Institutes of Health, National Eye Institute (nos. F30 EY025510 to E.Y.W., R01 EY026927 to A.S.T. and W.J.M., and T32 EY00252037 and T32 EY07001 to A.S.T.) and National Institute of Mental Health (nos. F30 F30MH088228 to R.J.C.). We thank F. Sinz for helpful discussion and suggestions on the DNN fitting to likelihood functions. We also thank T. Shinn for assistance in the behavioral training of the monkeys and experimental data collection.

Author contributions

All authors designed the experiments and developed the theoretical framework. R.J.C. programmed the experiment. 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. W.J.M. and A.S.T. supervised all stages of the project.

Competing interests

E.Y.W. and A.S.T. hold equity ownership in Vathes LLC, which provides development and consulting for the open source software (DataJoint) used to develop and operate a data analysis pipeline for this publication.

Additional information

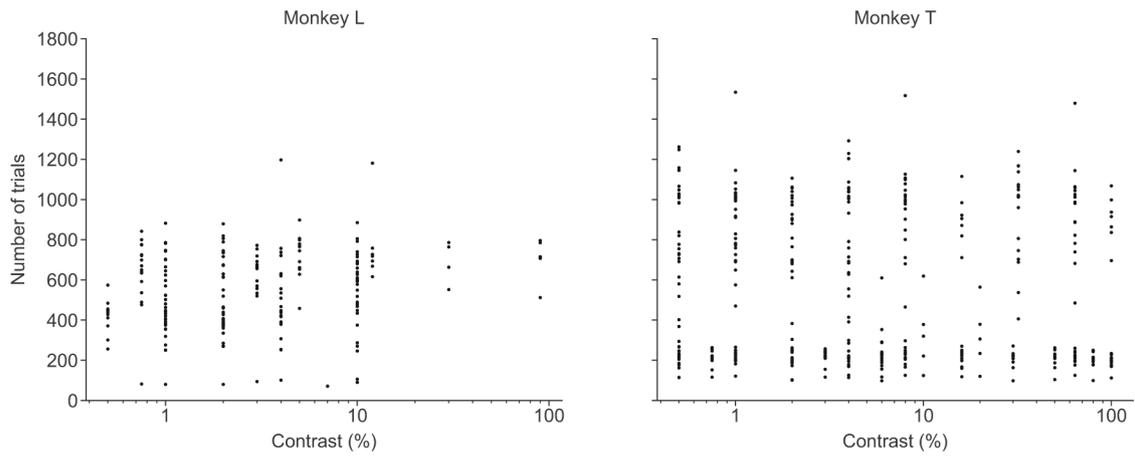
Extended data is available for this paper at <https://doi.org/10.1038/s41593-019-0554-5>.

Supplementary information is available for this paper at <https://doi.org/10.1038/s41593-019-0554-5>.

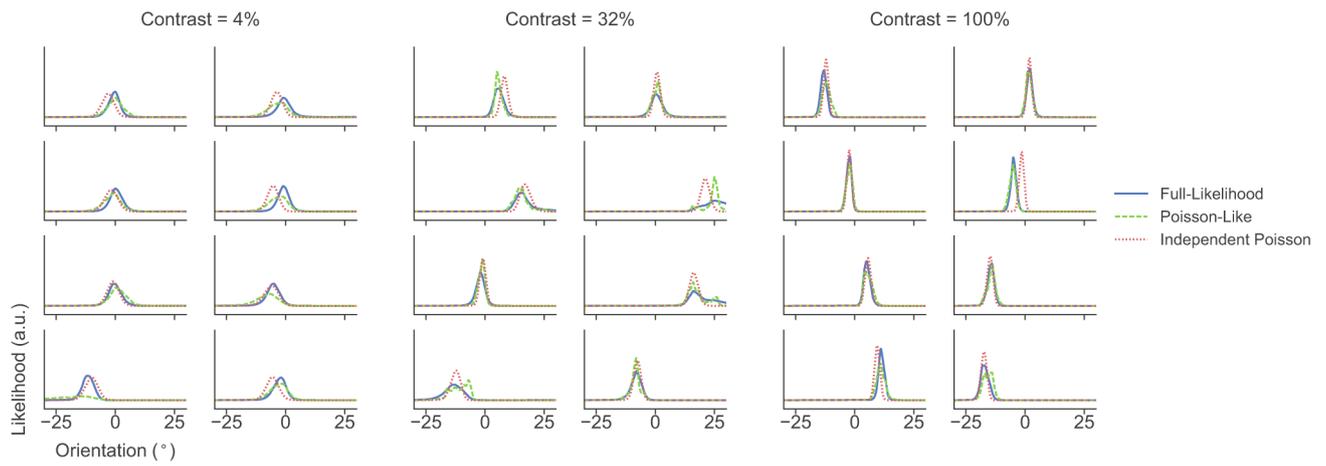
Correspondence and requests for materials should be addressed to E.Y.W., W.J.M. or A.S.T.

Peer review information *Nature Neuroscience* thanks Jan Drugowitsch and Robbe Goris for their contribution to the peer review of this work.

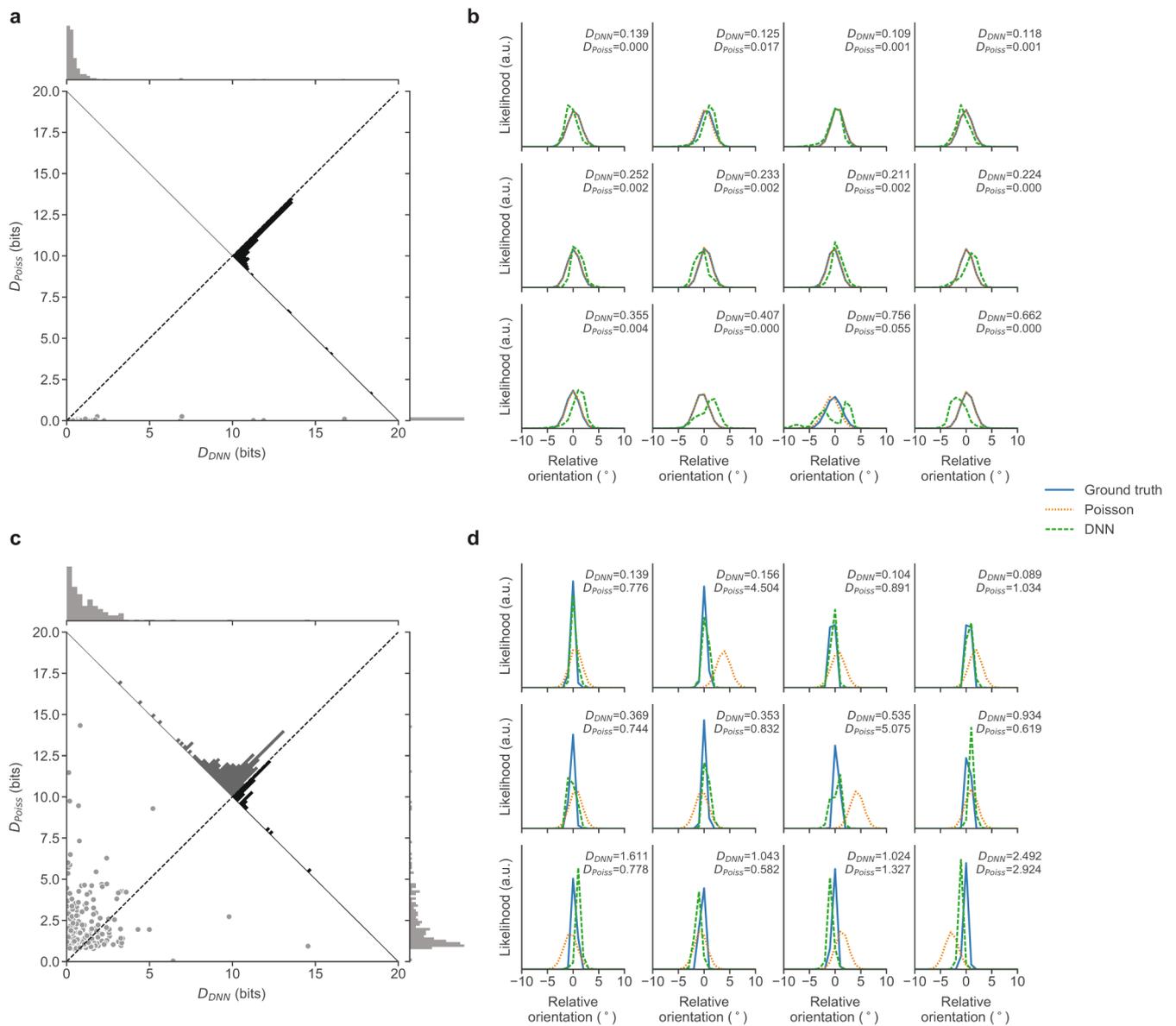
Reprints and permissions information is available at www.nature.com/reprints.



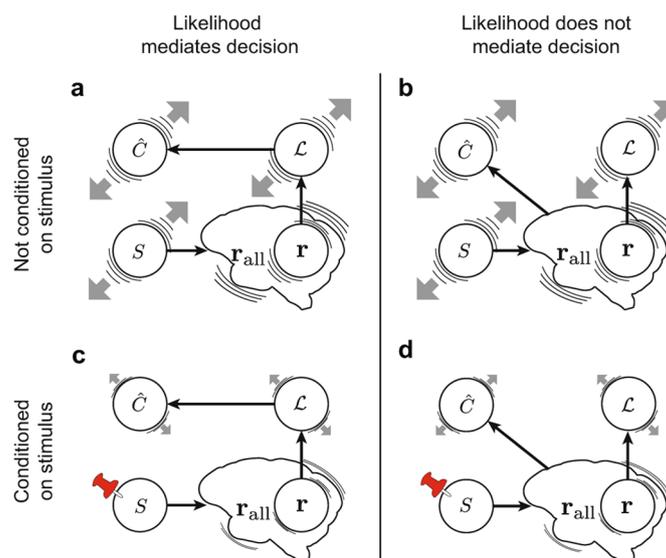
Extended Data Fig. 1 | Number of trials per contrast-session. Each point corresponds to a single contrast-session, depicting the number of trials performed at the particular contrast.



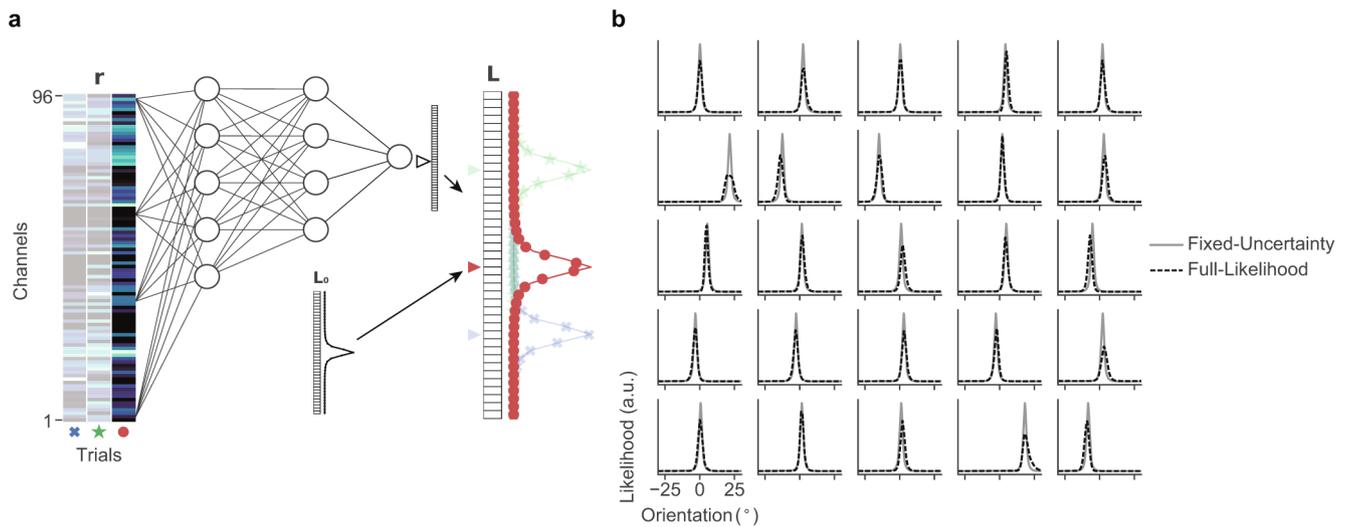
Extended Data Fig. 2 | Example decoded likelihood functions. Example decoded likelihood functions under Full-Likelihood, Poisson-like and Independent-Poisson based decoders are shown for randomly selected trials from three distinct contrast-sessions from Monkey T.



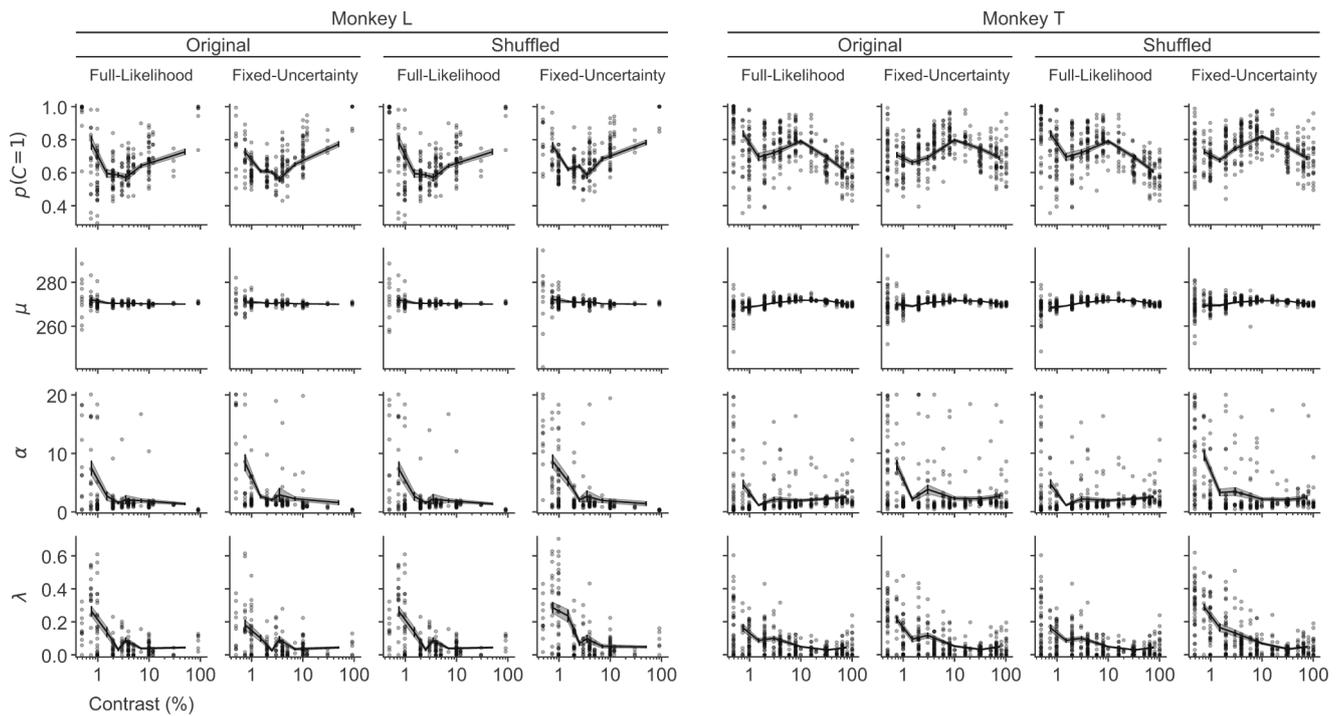
Extended Data Fig. 3 | Performance of the likelihood functions decoded by DNN-based decoders. a, b, Results on independent Poisson population responses. **a**, KL divergence between the ground truth likelihood function and likelihood function decoded with: a trained DNN D_{DNN} vs. independent Poisson distribution assumption D_{Pois} . Each point is a single trial in the test set. The distributions of D_{DNN} and D_{Pois} are shown at the top and right margins, respectively. The distribution of pair-wise difference between D_{DNN} and D_{Pois} 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_{DNN} (rows). **c, d**, Same as in **a, b** but for simulated population responses with correlated Gaussian distribution where the variance is scaled by the mean.



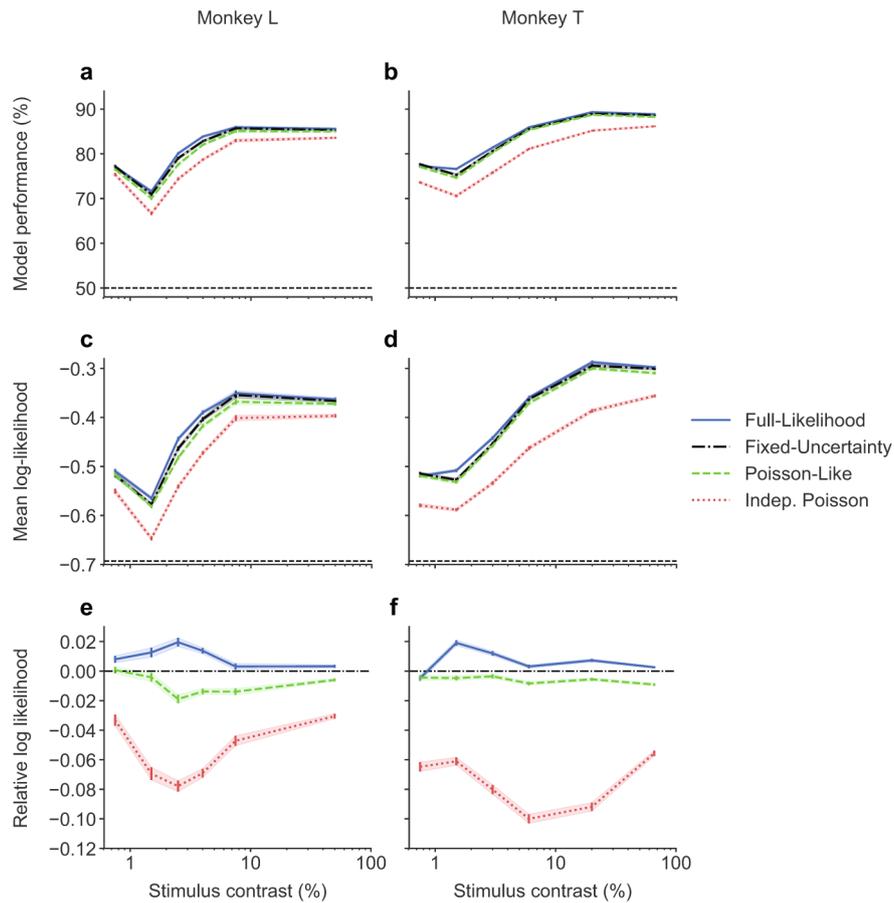
Extended Data Fig. 4 | Alternative relationships between the likelihood function and the decision. Possible relationships between variables in the model are indicated by black arrows. We consider two scenarios: **a, c** the likelihood function \mathcal{L} mediates the decision \hat{C} , **b, d** the 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 (red push pins), 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_{all} represents responses of all recorded and unrecorded neurons.



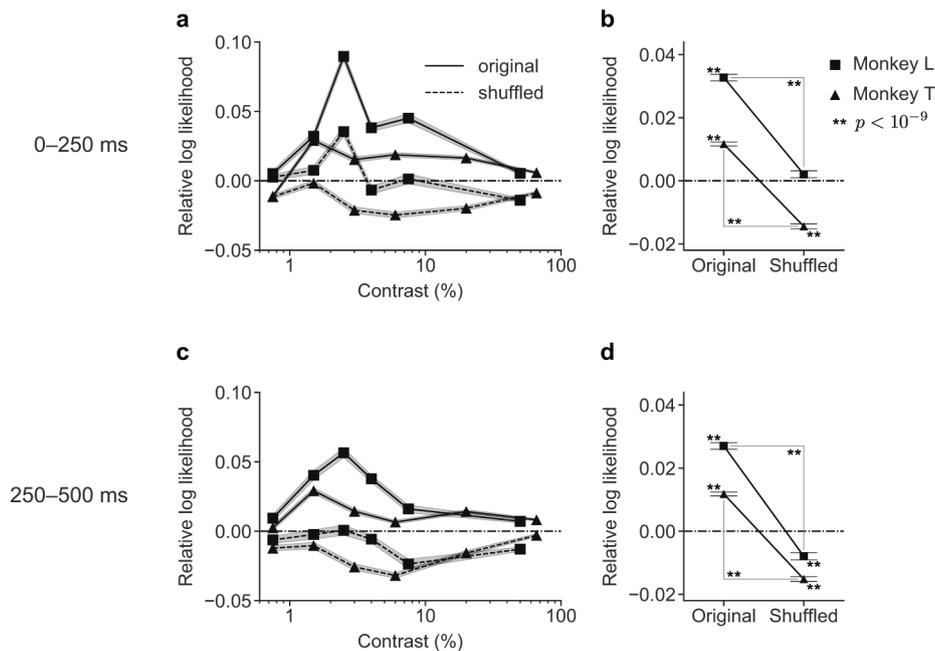
Extended Data Fig. 5 | Fixed-Uncertainty decoder. **a**, A schematic of a DNN for the Fixed-Uncertainty decoder mapping \mathbf{r} to the decoded likelihood function \mathbf{L} . For each contrast-session, the Fixed-Uncertainty decoder learns a single fixed-shape likelihood function \mathbf{L}_0 and a network that shifts \mathbf{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 trial. **b**, Example decoded likelihood functions from randomly selected trials from a single contrast-session for both the Fixed-Uncertainty decoder and the Full-Likelihood decoder.



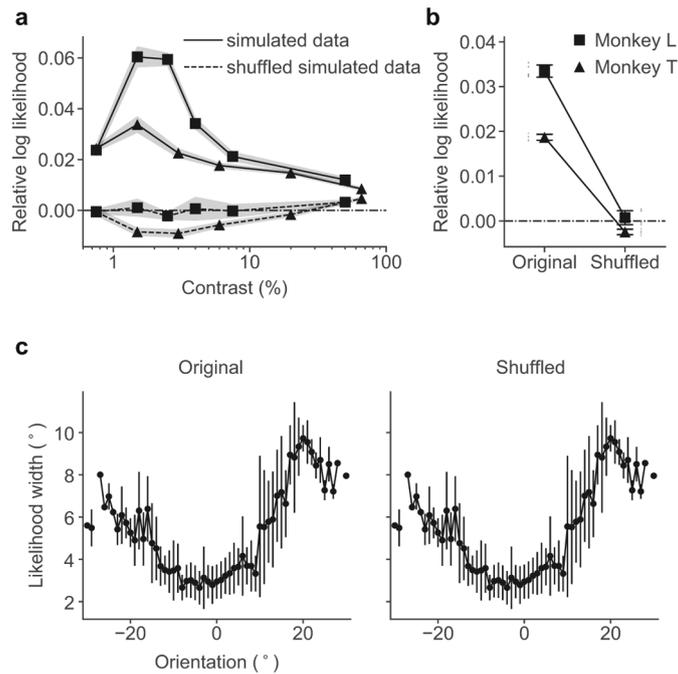
Extended Data Fig. 6 | Fitted Bayesian decision maker parameters. Each point corresponds to a single contrast-session, depicting the average fitted parameter value across 10 cross-validation training sets plotted against the contrast of the contrast-session. The solid line and error bars/shaded area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



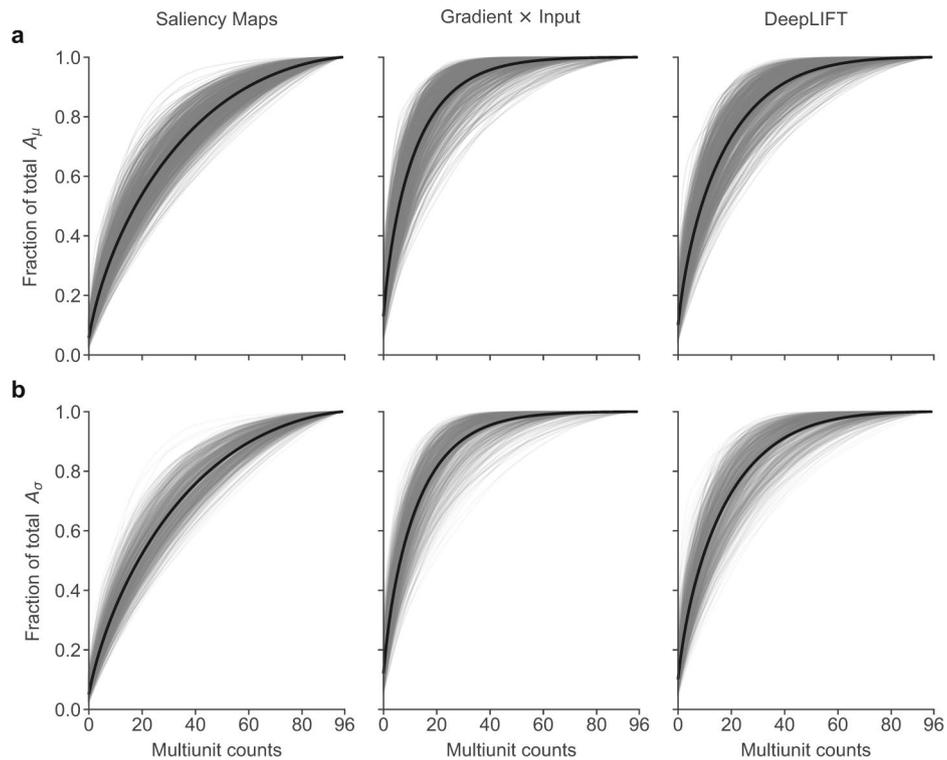
Extended Data Fig. 7 | Model performance on decision predictions. **a, b**, Model performance measured in proportions of trials correctly predicted by the model as a function of contrast for four decision models based on different likelihood decoders ($n=110,695$ and $n=192,630$ total trials across all contrasts for Monkey L and T, respectively). On each trial, the class decision that would maximize the posterior $P(\hat{C} | \mathbf{r})$ was chosen to yield a concrete classification prediction. **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**, black dashed lines indicate the performance at chance (50 % and $\ln(0.5)$, respectively). **e, f**, 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 ($n=110,695$ and $n=192,630$ total trials for Monkey L and T, respectively). Results are shown for the cross-validated datasets. All data points are the means and error bar/shaded area indicates the standard error of the mean.



Extended Data Fig. 8 | 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, ‘fh’) ($n=110,816$ and $n=192,962$ total trials for Monkey L and T) or **(c)** the second half (250–500 ms, ‘sh’) ($n=110,887$ and $n=192,980$ total trials for Monkey L and T) 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 trial log likelihood difference between the two models was statistically significant for both stimulus windows, and on both the original and the shuffled data for both monkeys (two tailed paired t -tests; Monkey L: $t_{\text{fh,original}}(110815) = 31.29$, $t_{\text{sh,original}}(110886) = 25.86$, $t_{\text{sh,shuffled}}(110886) = -6.98$; Monkey T: $t_{\text{fh,original}}(192961) = 18.48$, $t_{\text{fh,shuffled}}(192961) = -19.31$, $t_{\text{sh,original}}(192979) = 19.01$, $t_{\text{sh,shuffled}}(192979) = -20.17$; all with $p < 10^{-9}$), 0–250 ms for Monkey L ($t_{\text{fh,shuffled}}(110815) = 1.89$ with $P = 0.17$). The difference between the Full-Likelihood Model on the original and the shuffled data was significant for both monkeys for both stimulus windows (two tailed paired t -tests; Monkey L: $t_{\text{fh}}(110815) = 32.73$, $t_{\text{sh}}(110886) = 37.10$; Monkey T: $t_{\text{fh}}(192961) = 40.69$, $t_{\text{sh}}(192979) = 42.78$; all with $P < 10^{-9}$). All p values are Bonferroni corrected for the three comparisons. All data points are means, and error bar/shaded area indicate standard error of the means.



Extended Data Fig. 9 | Expected model performance on simulated data and observed effect of shuffling. **a, b**, Using the trained Full-Likelihood Model as the ground truth to simulate the behavior, the expected performances of the model on the simulated data was assessed. **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 indicates the standard deviation across the 5 simulation repetitions. For **b**, data points for individual simulation repetitions are depicted by gray icons next to the error bars. **c**, The dependence of the width of the likelihood function σ_l on the stimulus orientation is depicted for an example contrast-session (Monkey T, 8 % contrast, $n=1,126$ trials) on the original and the shuffled data. The shuffling procedure preserves the relationship between the average likelihood width and the stimulus orientation as desired. All data points are means, and error bar indicates standard deviation across trials falling in the specific bin.



Extended Data Fig. 10 | Contributions of multi-units to the total attribution. a, For each contrast-session, the multi-units were ordered from the largest to the smallest attribution to the likelihood mean A_{μ} , and the cumulative attribution over the total of 96 multi-units were plotted (thin gray lines, $n=545$ total contrast-sessions from Monkey L and T). The average cumulative attribution over all contrast-sessions are depicted by the thick black lines. The results are shown for each attribution method separately. **b**, Same as in **a**, but for attribution to the likelihood standard deviation A_{σ} .

Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Research policies, see [Authors & Referees](#) and the [Editorial Policy Checklist](#).

Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

n/a Confirmed

- The exact sample size (n) for each experimental group/condition, given as a discrete number and unit of measurement
- A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- The statistical test(s) used AND whether they are one- or two-sided
Only common tests should be described solely by name; describe more complex techniques in the Methods section.
- A description of all covariates tested
- A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
- A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
- For null hypothesis testing, the test statistic (e.g. F , t , r) with confidence intervals, effect sizes, degrees of freedom and P value noted
Give P values as exact values whenever suitable.
- For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
- For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
- Estimates of effect sizes (e.g. Cohen's d , Pearson's r), indicating how they were calculated

Our web collection on [statistics for biologists](#) contains articles on many of the points above.

Software and code

Policy information about [availability of computer code](#)

Data collection

Experimental data were collected using the following tools/software: LabVIEW 2017, MATLAB R2017a, and Psychtoolbox-3.0.12. Custom LabVIEW and MATLAB code used for stimulus presentation, electrophysiology data collection and data processing are made publicly available at the Tolias lab GitHub organization page <https://github.com/atlab>.

Data analysis

All analysis were performed using custom software developed using the following tools and softwares: MATLAB (R2017a), Python (3.6), DataJoint (0.11.1), PyTorch (0.4.1), NumPy (1.16.4), SciPy (1.3.0), Docker (18.09.7), matplotlib (3.0.3), seaborn (0.9.0), pandas (0.24.2), and Jupyter (1.0.0). Specifically, spike detection was performed using custom spike detection routine written in MATLAB and the code is available at <https://github.com/atlab/spikedetection>. All code used for modeling and training the deep neural networks as well as for figure generation will be viewable and available for download from https://github.com/eywalker/v1_likelihood. All other code used for analysis including data selection and decision model fitting can be found at https://github.com/eywalker/v1_project.

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors/reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Research [guidelines for submitting code & software](#) for further information.

Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A list of figures that have associated raw data
- A description of any restrictions on data availability

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

Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

Life sciences Behavioural & social sciences Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see [nature.com/documents/nr-reporting-summary-flat.pdf](https://www.nature.com/documents/nr-reporting-summary-flat.pdf)

Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	The sample size of two monkeys has been established following the standard practice in the non-human primate research and to ensure independent replication of results observed from the first monkey in the second monkey.
Data exclusions	Some recorded trials were excluded from the analysis if there was an excessive noise in the recorded signals. Details of the exclusion criteria can be found in the Methods, and the criteria was established prior to testing any models.
Replication	Every findings described in this study on one monkey has been successfully replicated on the second monkey that was independently trained, recorded and analyzed.
Randomization	There was no randomization performed as the study does not involve multiple study groups, and all analyses were performed in identical fashion on all subjects.
Blinding	There was no blinding performed as the study does not involve multiple study groups, and all analyses were performed in identical fashion on all subjects.

Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

Materials & experimental systems

n/a	Involvement in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> Antibodies
<input checked="" type="checkbox"/>	<input type="checkbox"/> Eukaryotic cell lines
<input checked="" type="checkbox"/>	<input type="checkbox"/> Palaeontology
<input type="checkbox"/>	<input checked="" type="checkbox"/> Animals and other organisms
<input checked="" type="checkbox"/>	<input type="checkbox"/> Human research participants
<input checked="" type="checkbox"/>	<input type="checkbox"/> Clinical data

Methods

n/a	Involvement in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> ChIP-seq
<input checked="" type="checkbox"/>	<input type="checkbox"/> Flow cytometry
<input checked="" type="checkbox"/>	<input type="checkbox"/> MRI-based neuroimaging

Animals and other organisms

Policy information about [studies involving animals](#); [ARRIVE guidelines](#) recommended for reporting animal research

Laboratory animals	This study involved two healthy, male rhesus macaque (<i>Macaca mulatta</i>) monkeys aged 10 and 7 years.
Wild animals	This study did not involve wild animals.
Field-collected samples	This study did not involve samples collected from the field.
Ethics oversight	All procedures were approved by the Institutional Animal Care and Use Committee (IACUC) of Baylor College of Medicine.

Note that full information on the approval of the study protocol must also be provided in the manuscript.