

# Population Code Dynamics in Categorical Perception

Chihiro I. Tajima<sup>a,\*</sup>, Satoshiro Tajima<sup>b,\*</sup>, Kowa Koida<sup>c</sup>, Hidehiko Komatsu<sup>d</sup>, Kazuyuki Aihara<sup>e,a</sup>, and Hideyuki Suzuki<sup>a</sup>

a. Graduate School of Information Science and Technology, the University of Tokyo. 7-3-1 Hongo, Bunkyo, Tokyo 113-8656, Japan.

b. Department of Basic Neuroscience, University of Geneva. CMU, 1 rue Michel Servet, 1211 Genève, Switzerland.

c. EIIRIS, Toyohashi University of Technology. 1-1 Hibarigaoka, Tempaku, Toyohashi, Aichi, 441-8580, Japan.

d. National Institute for Physiological Sciences. 38 Nishigonaka Myodaiji, Okazaki, Aichi, 444-8585, Japan.

e. Institute of Industrial Science, the University of Tokyo. 4-6-1 Komaba, Meguro, Tokyo 153-8505, Japan.

\*These authors contributed equally to this work.

Corresponding authors: CIT: [chihi.at.heart@gmail.com](mailto:chihi.at.heart@gmail.com); ST: [satoshiro.tajima@gmail.com](mailto:satoshiro.tajima@gmail.com), Department of Neuroscience, University of Geneva. CMU, 1 rue Michel Servet, 1211 Genève, Switzerland.

## Supplementary Information

### Contents

Supplementary Note .....	2
A. Neural network model that approximates dynamical categorical inference .....	2
Step1: Exact statistical inference .....	2
Step 2: Approximated statistical inference.....	3
Step 3: Implementation by recursive population code.....	4
B. Electrophysiological recording and data analysis .....	11
References .....	13
Supplementary Figures.....	14
Supplementary Figure S1. ....	14
Supplementary Figure S2. ....	15
Supplementary Figure S3. ....	16

# 1 Supplementary Note

## 2 A. Neural network model that approximates dynamical categorical inference

3 We formulate the categorical perception as an online inference problem, and we propose a neural network model that works as an  
4 inference algorithm approximately. While the problem formulation and the model themselves are presented in **Results** in main  
5 text, we explain in this Supplementary Note the reason why the model works as a hierarchical Bayesian inference algorithm;  
6 namely, we provide technical details as to (i) how the exact inference algorithm is derived from the problem formulation, (ii) how  
7 the approximated algorithm is derived, and (iii) how the approximated algorithm is implemented in a biologically plausible neural  
8 network.

9 We begin by deriving the optimal statistical inference of the stimulus category, where the category and the sensory evidence  
10 may change across time. We consider the task that a higher visual area infers stimulus value  $\theta_t$  and its category  $c_t$  at time  $t$   
11 from the history of bottom-up signals  $D_t \equiv \{\mathbf{r}_t, \dots, \mathbf{r}_0\}$  transmitted from the early visual system until time  $t$ , where there are a  
12 given number ( $n$ ) of categories  $c_t \in \{1, \dots, n\}$ . To solve this problem, we assume that the response activities  $\mathbf{r}_t$  of sensory  
13 neurons are received from the early visual system. The temporal evolution of the stimulus value  $\theta_t$  is dependent on its category  
14  $c_t$  and follows the graphical model shown in **Fig 1a**.

### 15 *Step1: Exact statistical inference*

16 The posterior probability of stimulus (e.g., hue of color)  $\theta_t$  is expressed as follows:

$$P[\theta_t|D_t] = P[\theta_t|\mathbf{r}_t, D_{t-1}] \propto P[\mathbf{r}_t|\theta_t] P[\theta_t|D_{t-1}]. \quad (S1)$$

17 Considering the marginalization and dependencies among the variables, we have

$$P[\theta_t|D_t] \propto P[\mathbf{r}_t|\theta_t] \sum_{c_t} \sum_{\theta_{t-1}} P[\theta_t|\theta_{t-1}, c_t] P[\theta_{t-1}|D_{t-1}] \sum_{c_{t-1}} P[c_t|c_{t-1}] P[c_{t-1}|D_{t-1}]. \quad (S2)$$

18 This representation of the posterior contains deeply nested summations. Given the potential computational cost and restrictions on  
19 the nervous system, it is not likely that the biological system implement the calculation shown in Eq. (S2) *in situ*<sup>1</sup>. In the  
20 subsequent sections, we derive an approximated representation of the posterior.

1 **Step 2: Approximated statistical inference**

2 We propose a reduced estimation scheme that approximates the exact inference shown above. We assume that a neural  
3 population receive sensory input signals containing stimulus information and noise. The purpose of the model is to estimate the  
4 true stimulus identity without noise. We introduce the following three approximations:

5 (a1) *The category slowly changes compared to the stimulus:*

$$P[c_t | c_{t-1}] \simeq \delta_{c_t, c_{t-1}}. \quad (S3)$$

6 In natural environment, it is reasonable to assume that the stimulus can change, but the category rarely changes over time except  
7 for drastic changes as explained above. For example, appearance of object surface can be determined by the surface reflectance  
8 property and the lighting condition. Although the lighting condition (e.g., the surface angle) can alter in a short time, the surface  
9 property (e.g., reflectance spectra) is often relatively constant. Therefore, it is relatively rare to observe a drastic appearance change  
10 straddling a categorical boundary (e.g., such that a red object suddenly becomes blue) even though the precise appearance of  
11 surface may change rapidly.

12 (a2) *Omitting the uncertainty of the previous category estimate:*

$$P[c_{t-1} | D_{t-1}] \simeq \delta_{c_{t-1}, \hat{c}_{t-1}}, \quad (S4)$$

13 where  $\hat{c}_{t-1}$  is the categorical estimate at time  $t - 1$ . As described later, in the simulation, we provided the categorical estimate  
14 at each time point by maximizing the likelihood of sensory input.

15 (a3) *Variable separation:* The posterior probabilities of an estimated stimulus, conditioned by  $\theta_{t-1}$  and  $c_t$ , were assumed to be  
16 approximated as

$$P[\theta_t | \theta_{t-1}, c_t] \propto P[\theta_t | \theta_{t-1}] P[\theta_t | c_t]. \quad (S5)$$

17 This approximation is valid in many natural situations, for example, where  $P[\theta_t | \theta_{t-1}, c_t]$ ,  $P[\theta_t | \theta_{t-1}]$ , and  $P[\theta_t | c_t]$  are all  
18 Gaussian-like unimodal distributions and  $P[\theta_t | \theta_{t-1}, c_t]$  is maximized at a point between the peaks of  $P[\theta_t | \theta_{t-1}]$  and  
19  $P[\theta_t | c_t]$ .

20 Although these three assumptions are natural for realistic environments, it should be noted that they do not always hold,

1 especially in laboratory experiments. For example, in an experiment where the experimenter can use a completely unstructured  
 2 stimulus sequence, the slower dynamics of category compared to that of hue, Eq. (S3), could deviate from the reality. In this sense,  
 3 the inference based on this assumption becomes sub-optimal in such an artificial situation. Nevertheless, it is possible that the  
 4 brain, which could be adapted to the natural environment, still uses this sub-optimal inference strategy during the laboratory  
 5 experiment. Similarly, the assumption described by Eq. (S5) is not always valid, for example, when  $P[\theta_t|\theta_{t-1}, c_t]$  is not  
 6 maximized between the peaks of  $P[\theta_t|\theta_{t-1}]$  and  $P[\theta_t|c_t]$ . This could occur when the experimenter artificially introduces a  
 7 complex structure of category such that  $P[\theta_t|c_t]$  is a multimodal function of hue  $\theta_t$ . The inference under such the complex  
 8 categorical structure is beyond the scope of the present study although it is an interesting direction of future extension of the  
 9 current model.

10 Applying the above approximations (a1)–(a3) to Eq. (S2), we have:

$$P[\theta_t|D_t] \propto P[\mathbf{r}_t|\theta_t] \sum_{\theta_{t-1}} P[\theta_t|\theta_{t-1}] P[\theta_{t-1}|D_{t-1}] P[\theta_t|c_t = \hat{c}_{t-1}], \quad (\text{S6})$$

11 or, equivalently,

$$\ln P[\theta_t|D_t] \simeq \ln P[\mathbf{r}_t|\theta_t] + \ln \sum_{\theta_{t-1}} P[\theta_t|\theta_{t-1}] P[\theta_{t-1}|D_{t-1}] + \ln P[\theta_t|c_t = \hat{c}_{t-1}] + \text{const.} \quad (\text{S7})$$

12 This corresponds to Eq. (3) in **Results** in main text.

### 13 *Step 3: Implementation by recursive population code*

14 Next, we explain how the above online statistical inference can be achieved with a recurrent neural network (the parameters used  
 15 in the simulation are summarized in **Table S1**). Let  $r_t^i$  denote each bottom-up, stimulus-evoked spike signal (hereafter, we refer  
 16 to it as *sensory input*), received by the  $i$ th hue-selective neuron at time point  $t$ . We assume that the sensory inputs to individual  
 17 neurons,  $\mathbf{r}_t = (r_t^1, \dots, r_t^N)$ , are generated by an independent Poisson process for a given stimulus value,  $\theta_t$  as follows:

$$P[\mathbf{r}_t|\theta_t] = \prod_i \frac{f(\theta_t - \varphi_i)^{r_t^i}}{r_t^i!} e^{-f(\theta_t - \varphi_i)}, \quad (\text{S8})$$

18 where  $f(\theta_t - \varphi_i)$  is the expected spike count between steps  $t - 1$  and  $t$  in the  $i$ th hue-selective neuron, and  $\varphi_i$  denotes  
 19 its preferred stimulus. (Here we assume the independence of spike-count variability in  $\mathbf{r}_t$  given the input to each neuron, but this  
 20 should be discriminated from the apparent noise correlation that may include fluctuations in shared network input to each neuron.  
 21 The assumption of independent Poisson process does not rule out such the correlated variability reflecting the correlation in  
 22 network input to neurons; rather, it focuses on the variability within each neuron after eliminating the global fluctuation that



1 originates from shared signals<sup>2</sup>.)

2 If we assume that the preferred stimulus is uniformly distributed across the neural population, the log-likelihood of the given  
3 sensory input is simplified as:

$$\ln P[\mathbf{r}_t | \theta_t] = \sum_i r_t^i \ln f(\theta_t - \varphi_i) + \text{const.} \quad (\text{S9})$$

4 Here, the constant term is independent of  $\theta_t$ . We used the fact that the summation  $\sum_i f(\theta_t - \varphi_i)$  does not depend on  $\theta_t$  due  
5 to the uniformity when the total number of neurons is sufficiently large. (Note that the concept of uniformity depends on the  
6 choice of the stimulus space; we assume a stimulus space that approximately satisfies this condition. In addition, the uniformity  
7 assumption about neuronal selectivity often does not hold for intensity dimensions, such as stimulus contrast.) Notably, Eq. (S9)  
8 implies a mapping relationship from a population input  $\mathbf{r}$  to the logarithm of the distribution function of stimulus  $\theta$ , through the  
9 linear summation of kernels,  $\ln f(\theta_t - \varphi_i)$ <sup>3-7</sup>. Assuming a generic nonlinear function  $f$  and a sufficient number of neurons  
10 that prefer different stimulus values, arbitrary functions of  $\theta$  can be mapped to the space of the population input  $\mathbf{r}$ . Specific  
11 constraints on  $f$  lead to the restriction on the image of this map, but the image includes a family of functions practically sufficient  
12 to encode posterior distributions of the stimulus under the biologically plausible assumptions explained below.

13 To be specific, we introduce parametrized models of the neural tuning function and the probability distributions. First, we  
14 assume that the neural tuning function  $f$  is a von Mises function (circular Gaussian),  $f(x) = \exp(\kappa \cos x)$ , with a tuning  
15 sharpness parameter  $\kappa$ . Second, we consider the situation in which the stimulus transition is described as  $P[\theta_t | \theta_{t-1}] \propto$   
16  $\exp(\sigma \cos(\theta_t - \theta_{t-1}))$  with a relatively large sharpness  $\sigma$ . Third, the top-down prior  $P[\theta_t | c_t = \hat{c}_{t-1}]$  is described by a  
17 von Mises distribution with an arbitrary sharpness  $\kappa^{\text{cat}}$  and a mode (or *focal stimulus*)  $\varphi_{\hat{c}_{t-1}}^{\text{cat}}$  that corresponds to the previously  
18 estimated category  $\hat{c}_{t-1}$ :

$$P[\theta_t | c_t = \hat{c}_{t-1}] \propto \exp(\kappa^{\text{cat}} \cos(\theta_t - \varphi_{\hat{c}_{t-1}}^{\text{cat}})). \quad (\text{S10})$$

19 At the very beginning of inference, it is reasonable to assume that the nervous system has no prior information about the  
20 stimulus, and that the activities of the hue-selective neurons  $\boldsymbol{\rho}_t = (\rho_t^1, \dots, \rho_t^N)$  at time  $t$  represent the likelihood function  
21 simply reflecting the sensory inputs:  $\rho_t^i = r_t^i$  ( $\forall i$ ). Our basic idea in the present study is that we can also represent the posterior

1 distribution of  $\theta$  by modifying the population input  $\mathbf{r}$  in Eq. (S9) to reflect prior distributions that depend on previous input  
 2 history and on top-down signals, corresponding to the second and the third terms in Eq. (S6); that is, there exist vectors  
 3  $\mathbf{a}_t = (a_t^1, \dots, a_t^N)$  and  $\mathbf{b}_t = (b_t^1, \dots, b_t^N)$  that satisfy

$$\ln \sum_{\theta_{t-1}} P[\theta_t | \theta_{t-1}] P[\theta_{t-1} | D_{t-1}] = \sum_i a_t^i \ln f(\theta_t - \varphi_i) + \text{const.}, \quad (\text{S11})$$

$$\ln P[\theta_t | c_t = \hat{c}_{t-1}] = \sum_i b_t^i \ln f(\theta_t - \varphi_i) + \text{const.} \quad (\text{S12})$$

4 Then, Eq. (S7) is rewritten as follows:

$$\ln P[\theta_t | D_t] = \sum_i \rho_t^i \ln f(\theta_t - \varphi_i) + \text{const.}, \quad (\text{S13})$$

$$\rho_t^i = r_t^i + a_t^i + b_t^i, \quad (\text{S14})$$

5 where  $\rho_t^i$  is the updated activity of the  $i$ th hue-selective neuron at time  $t$ . In this equation, the product of the distribution  
 6 functions  $P[\mathbf{r}_t | \theta_t] P[\theta_t | c_t = \hat{c}_{t-1}] \sum_{\theta_{t-1}} P[\theta_t | \theta_{t-1}] P[\theta_{t-1} | D_{t-1}]$  is substituted by the linear sum,  $r_t^i + a_t^i + b_t^i$ , where  
 7  $a_t^i$  and  $b_t^i$  are interpreted as bias inputs to the neuron  $i$ . Now, the problem is to derive the appropriate functional forms of  $a_t^i$   
 8 and  $b_t^i$ .

9 First, we show that,  $b_t^i$  is expressed in the following form:

$$b_t^i = \beta f^{\text{cat}}(\varphi_i - \varphi_{\hat{c}_{t-1}}^{\text{cat}}), \quad (\text{S15})$$

10 with an arbitrary even function  $f^{\text{cat}}$  if  $\beta$  satisfies the following condition:

$$\beta = \frac{\kappa^{\text{cat}} \cos(\theta_t - \varphi_{\hat{c}_{t-1}}^{\text{cat}})}{\kappa \sum_i f^{\text{cat}}(\varphi_i - \varphi_{\hat{c}_{t-1}}^{\text{cat}}) \cos(\theta_t - \varphi_i)} = \frac{\kappa^{\text{cat}}}{\kappa \tilde{f}^{\text{cat}}}, \quad (\text{S16})$$

11 where we defined  $\tilde{f}^{\text{cat}}$  as the cosine coefficient in the Fourier transforms of  $f^{\text{cat}}$ . To see this, using Eqs. (S10) and (S15) and  
 12 the definition of tuning function,  $f(x) = \exp(\kappa \cos x)$ , (S15) can be rewritten as

$$\kappa^{\text{cat}} \cos(\theta_t - \varphi_{\hat{c}_{t-1}}^{\text{cat}}) = \sum_i \beta f^{\text{cat}}(\varphi_i - \varphi_{\hat{c}_{t-1}}^{\text{cat}}) \kappa \cos(\theta_t - \varphi_i), \quad (\text{S17})$$

13 ignoring the constant term that is irrelevant to the hue estimation. This gives the first equality in Eq. (S16). To derive the second  
 14 equality in Eq. (S16), we can rewrite the summation in denominator of the middle term as  $\sum_i f^{\text{cat}}(\theta - \varphi_{\hat{c}_{t-1}}^{\text{cat}} - \psi) \cos \psi$  by

1 defining  $\psi = \theta - \varphi_i$ . From the uniformity of preferred-stimulus distribution across neurons, this summation is interpreted as the  
2 convolution between  $f^{\text{cat}}$  and  $\cos$ , and thus simply expressed as  $\tilde{f}^{\text{cat}} \cos(\theta - \varphi_{\hat{c}_{t-1}}^{\text{cat}})$  with the cosine coefficient  $\tilde{f}^{\text{cat}}$ ,  
3 using that  $f^{\text{cat}}$  is an even function. Now, the  $\theta$ -dependent term,  $\cos(\theta - \varphi_{\hat{c}_{t-1}}^{\text{cat}})$ , cancels out between the numerator and the  
4 denominator, providing the second equality in Eq. (S16). Therefore, we have a  $\theta$ -independent expression,  $\beta = \kappa^{\text{cat}} / \kappa \tilde{f}^{\text{cat}}$ ,  
5 and there exist  $\beta$  that satisfies this condition. Equation (S16) illustrates that the gain of the bias input should be determined as  
6 proportional to the certainty of the prior; e.g., when the top-down prior knowledge is certain, the sharpness parameter  $\kappa^{\text{cat}}$  of the  
7 distribution function  $P[\theta_t | \hat{c}_{t-1}]$  has a large value, and also the gain of top-down bias  $\beta$  should be large. In the simulation, we  
8 assumed  $\kappa^{\text{cat}} = \kappa$  for the simplicity although  $\kappa^{\text{cat}}$  is not necessarily equal to  $\kappa$  if the Eq. (S16) is satisfied.

9 Next, we derive  $a_t^i$ . Similarly to Eq. (S9), the log-posterior of the past stimulus is expressed (using the von Mises tuning  
10 property and the neural activity  $\rho_{t-1}^i$ ) as

$$\ln P[\theta_{t-1} | D_{t-1}] = \sum_i \rho_{t-1}^i \cos(\theta_{t-1} - \varphi_i) + \text{const.} = \tilde{\rho}_{t-1} \cos(\theta_{t-1} - \hat{\theta}_{t-1}) + \text{const.}, \quad (\text{S18})$$

11 with  $\hat{\theta}_{t-1}$ , the mode of the distribution, and  $\tilde{\rho}_{t-1}$ , the cosine coefficient in the Fourier transforms of activity as a function of  
12 preferred stimulus:  $\rho_{t-1}(\varphi_i) = \rho_{t-1}^i$ . Eq. (S18) and  $P[\theta_t | \theta_{t-1}] \propto \exp(\sigma \cos(\theta_t - \theta_{t-1}))$  yield  
13  $\sum_{\theta_{t-1}} P[\theta_t | \theta_{t-1}] P[\theta_{t-1} | D_{t-1}] \propto \sum_{\theta_{t-1}} \exp(\sigma \cos(\theta_t - \theta_{t-1})) \exp(\tilde{\rho}_{t-1} \cos(\theta_{t-1} - \hat{\theta}_{t-1}))$ , where the convolution  
14 between two von Mises functions in the left-hand side is approximated, for sufficiently large  $\sigma$  and  $\tilde{\rho}_{t-1}$ , as

$$\begin{aligned} \ln \sum_{\theta_{t-1}} P[\theta_t | \theta_{t-1}] P[\theta_{t-1} | D_{t-1}] &\approx \frac{\sigma \tilde{\rho}_{t-1}}{\tilde{\rho}_{t-1} + \sigma} \cos(\theta_t - \hat{\theta}_{t-1}) + \text{const.} \\ &= \frac{\sigma}{\tilde{\rho}_{t-1} + \sigma} \sum_i \rho_{t-1}^i \ln f(\theta_t - \varphi_i) + \text{const.} \end{aligned} \quad (\text{S19})$$

15 (To show the approximated equality, for example, we can consider a Gaussian approximation of von Mises function, such as  
16  $\exp(\sigma \cos(\theta_t - \theta_{t-1})) \propto \exp(-\sigma(\theta_t - \theta_{t-1})^2/2)$ , and a convolution of two Gaussian functions  
17  $\sum_{\theta_{t-1}} \exp(-\sigma(\theta_t - \theta_{t-1})^2/2) \exp(-\tilde{\rho}_{t-1}(\theta_{t-1} - \hat{\theta}_{t-1})^2/2) \propto \exp(-\frac{\sigma \tilde{\rho}_{t-1}}{\tilde{\rho}_{t-1} + \sigma}(\theta_t - \hat{\theta}_{t-1})^2/2)$ .) Comparing it to  
18 Eq. (S11), we have

$$a_t^i = \alpha_{t-1} \rho_{t-1}^i, \quad (\text{S20})$$

1 where  $\alpha_{t-1} = \sigma / (\tilde{\rho}_{t-1} + \sigma)$  can be implemented by the divisive gain-control mechanism in biological systems<sup>8</sup>. This term is  
 2 further approximated by a constant value if the ratio  $\tilde{\rho}_{t-1} / \sigma$  is within a relatively narrow range; i.e., the uncertainty of  
 3 momentary posterior and the magnitude of stimulus fluctuation are at roughly the same order, which would be reasonable in  
 4 practical situations.

5 Together,

$$\rho_t^i = r_t^i + \alpha \rho_{t-1}^i + \beta f^{\text{cat}}(\varphi_i - \varphi_{\hat{c}_{t-1}}^{\text{cat}}). \quad (\text{S21})$$

6 This corresponds to Eq. (4) in the **Results** in main text.. In the present simulations, we computed the quantities appearing in the  
 7 **Results** section—such as the stimulus discrimination threshold (via the ideal observer analysis, **Methods**) and the mean peak of  
 8 population activity—based on the spike statistics (mean, variance etc.) determined by the Poisson process, where each neuron’s  
 9 mean spike count is  $\rho_t^i$ . Note that, in the current problem setup, the expected value of  $\mathbf{r}$  is the sufficient statistic that describes  
 10 Poisson distribution over spike count variability of each neuron. In our simulation, we set the weights  $(\alpha, \beta) = (0.5, 0.2)$  to  
 11 roughly fit the data for single neuron response modulation (**Fig. 2**); moderate changes in these and other variables (including the  
 12 width of color selectivity and anisotropy of top-down connectivity) did not affect the main findings of the present paper.

13 Finally, we describe how the category is estimated based on the neural population activity. The categorical estimate is given  
 14 by maximizing the posterior based on the history of sensory input:

$$\hat{c}_t = \underset{c_t}{\operatorname{argmax}} P[c_t | D_t] = \underset{c_t}{\operatorname{argmax}} \sum_{\theta_t} P[c_t | \theta_t, D_t] P[\theta_t | D_t], \quad (\text{S22})$$

15 which is implemented as a winner-take-all competition in the neural network model (**Methods**). From the assumption that  
 16  $P[c_t | \theta_t, D_t]$  is approximated by a von Mises distribution over  $\theta_t$  that is maximized around  $\theta_t = \varphi_{c_t}^{\text{cat}}$ , the maximization  
 17 procedure is well approximated by:

$$\hat{c}_t = \underset{c_t}{\operatorname{argmax}} \ln P[\theta_t = \varphi_{c_t}^{\text{cat}} | D_t] = \underset{c_t}{\operatorname{argmax}} \sum_i \rho_t^i \ln f(\varphi_i - \varphi_{c_t}^{\text{cat}}). \quad (\text{S23})$$

18 Note that this approximation works not only for a von Mises distribution but for a general bell-shaped function  $P[\theta_t | c_t] =$

1  $g(\theta_t - \varphi_{c_t}^{\text{cat}})$  that peaks at the focal stimulus  $\varphi_{c_t}^{\text{cat}}$ , which provides a reasonable description of the categorical generative  
2 model for color or other stimulus features in general (e.g., orientation, motion, or facial expression). Equation (S23) demonstrates  
3 that the category estimate is simply provided by reading out the population activity of hue-selective neurons ( $\rho_t^i$ ) on each time step,  
4 as we implemented in the neural network mode (**Fig. 1d**). This is the case because the instantaneous population activity of  
5 hue-selective neurons reflects the history of previous category estimates: as we have described, the population activity of  
6 hue-selective neurons is modulated by the previous estimate of category as well as the previous activities of themselves, to  
7 represent the posterior distribution of hue ( $P[\theta_t | D_t]$ ) based on the history of sensory input ( $D_t$ ) up to that moment.

8

9

**Table S1.** Summary of model parameters used in the simulation.

Description	Notation	Value	Comment
Sharpness of tuning in hue-selective neurons	$\kappa$	3	
Sharpness of top-down connectivity from category-selective neurons to hue-selective neurons	$\kappa^{\text{cat}}$	3	
Strength of lateral interaction (relative to sensory input)	$\alpha$	0.5	
Strength of top-down interaction (relative to sensory input)	$\beta$	0.2	Set $\beta=0$ to simulate ‘Without top-down’ condition in Fig. 2.
Number of hue-selective neurons	$N$	300	
Number of category-selective neurons	$n$	3	
Preferred hue of hue-selective neuron $i$ ( $i \in \{1, \dots, N\}$ )	$\varphi_i$	$-\pi + 2\pi i/N$	Uniformly distributed in angle from $-\pi$ to $\pi$ (radians).
Preferred hue of category-selective neuron $i$ ( $i \in \{1, \dots, n\}$ )	$\varphi_i^{\text{cat}}$	$-\pi + 2\pi i/n$	

## 1 **B. Electrophysiological recording and data analysis**

2 Details of the surgical and recording procedures have been previously published<sup>9</sup>. Two female monkeys (*Macaca fuscata*) were  
3 used for the experiments. The monkeys were trained in a categorization task, a discrimination task, and a simple fixation task. In  
4 all three tasks, 11 sample colors were presented in a pseudorandom order. The monkey was required to maintain fixation within  
5 the trial, except for the saccade response. The sample color stimulus was presented for 500 ms. There were eleven sample colors  
6 that ranged from red [color 1, (x, y) = (0.631, 0.343)] to green [color 11, (x, y) = (0.286, 0.603)] with equal spaces on the  
7 International Commission on Illumination (CIE) xy chromaticity diagram.

8 In the categorization task, the monkey reported whether the sample color was reddish (sample colors 1–4) or greenish  
9 (sample colors 8–11) by saccade, and was rewarded for correct responses. For the intermediate colors (sample colors 5–7), the  
10 monkey was rewarded randomly regardless of its behavioral response. In the discrimination task, the monkey reported which test  
11 color was the same as the reference color by saccade. The two choice colors were three steps apart along the 11 sample colors: the  
12 eight choice-color pairs included colors 1–4, 2–5, 3–6, 4–7, 5–8, 6–9, 7–10, and 8–11. This color interval was chosen so as to  
13 yield a modest performance (about 80–90% correct).

14 Neuronal activity was recorded from the anterior part of the IT cortex, which is a region where color-selective neurons are  
15 concentrated. To record single unit activities, microelectrodes were inserted, and the activities of single neurons were isolated by  
16 matching spike templates. We analyzed only data from correct trials. The visual response to a sample stimulus was computed as  
17 the firing rate between 50 and 550 ms after the sample onset. To determine the neuronal color preferences, we first averaged the  
18 firing rates during the above time range, and selected the stimulus color (either of the 11 sample stimulus) that evoked the  
19 maximum firing rate for each cell, as its preferred color. In the derivation of the population response distribution, the individual  
20 cell responses within each time bin were divided by the time-averaged response to its preferred color during the fixation task.  
21 Then, the responses for cells that had the same color preference were respectively averaged. The peak loci of the population  
22 responses were obtained by fitting Gaussian functions with variable mean, variance, gain, and baseline level. The responses to the  
23 marginal colors (1 and 11 for the categorization task; 1, 2, and 11 for the fixation task) were excluded from the peak analysis  
24 because the peak estimates were not reliable for these stimuli. To test the statistical significance of the population activity  
25 modulation, we conducted 3-way analysis of variance (ANOVA; preferred stimulus×presented stimulus×time) on the difference

1 between the activities during categorization and discrimination tasks.

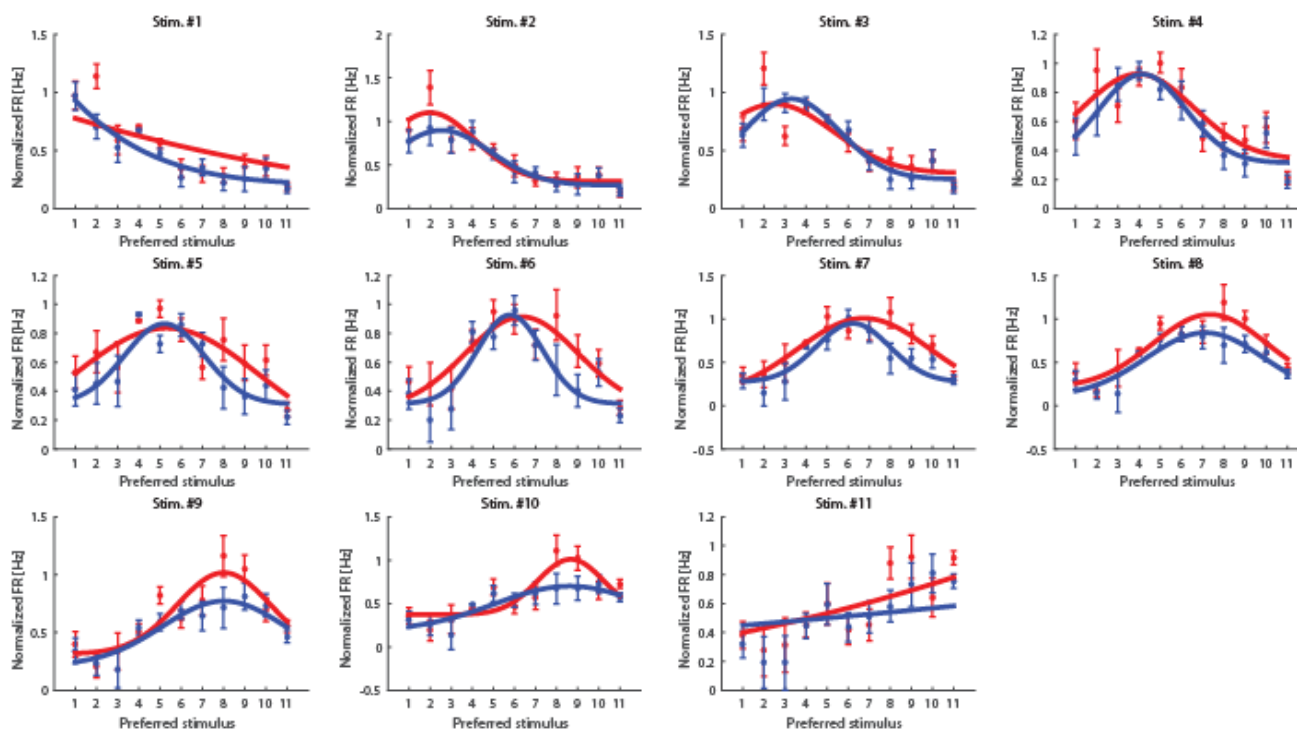
2 All procedures for animal care and experimentation were in accordance with the National Institutes of Health Guide for the  
3 Care and Use of Laboratory Animals and were approved by Institutional Animal Care and Use Committee of the National  
4 Institute of Natural Sciences.



## References

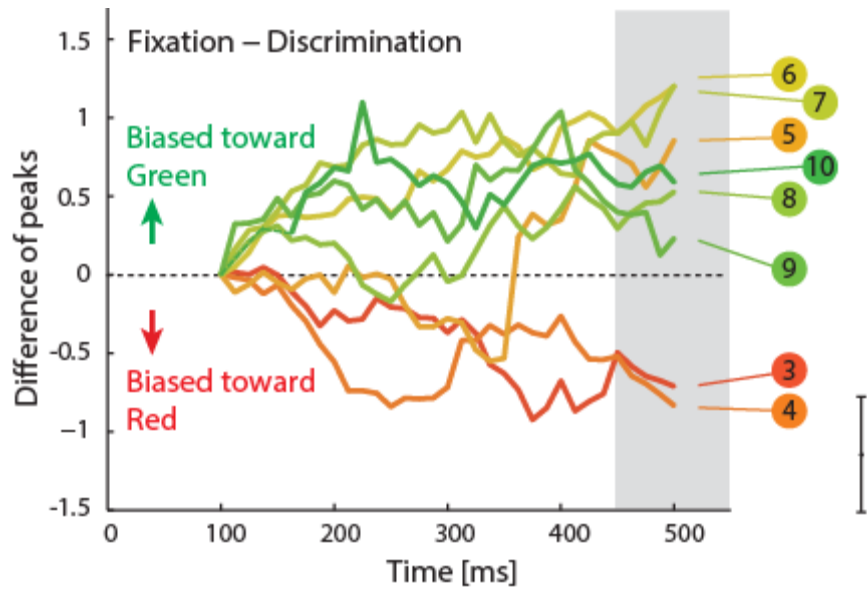
1. Yu, A. J. & Dayan, P. Acetylcholine in cortical inference. *Neural Networks* **15**, 719–30 (2002).
2. Goris, R. L. T., Movshon, J. A. & Simoncelli, E. P. Partitioning neuronal variability. *Nat. Neurosci.* **17**, 858–865 (2014).
3. Földiák, P. in *Computation and Neural Systems* (eds. Eeckman, F. H. & Bower, J. M.) **1992**, 55–60 (Norwell, MA: Kluwer Academic Publishers, 1993).
4. Sanger, T. Probability density estimation for the interpretation of neural population codes. *J. Neurophysiol.* **76**, 2799–2793 (1996).
5. Dayan, P. & Abbott, L. F. *Theoretical Neuroscience: Computational and Mathematical Modeling of Neural Systems*. (MIT Press., 2001).
6. Ma, W. J., Beck, J. M., Latham, P. E. & Pouget, A. Bayesian inference with probabilistic population codes. *Nat. Neurosci.* **9**, 1432–1438 (2006).
7. Jazayeri, M. & Movshon, J. A. Optimal representation of sensory information by neural populations. *Nat. Neurosci.* **9**, 690–696 (2006).
8. Beck, J. M., Latham, P. E. & Pouget, A. Marginalization in neural circuits with divisive normalization. *J. Neurosci.* **31**, 15310–15319 (2011).
9. Koida, K. & Komatsu, H. Effects of task demands on the responses of color-selective neurons in the inferior temporal cortex. *Nat. Neurosci.* **10**, 108–16 (2007).

## Supplementary Figures



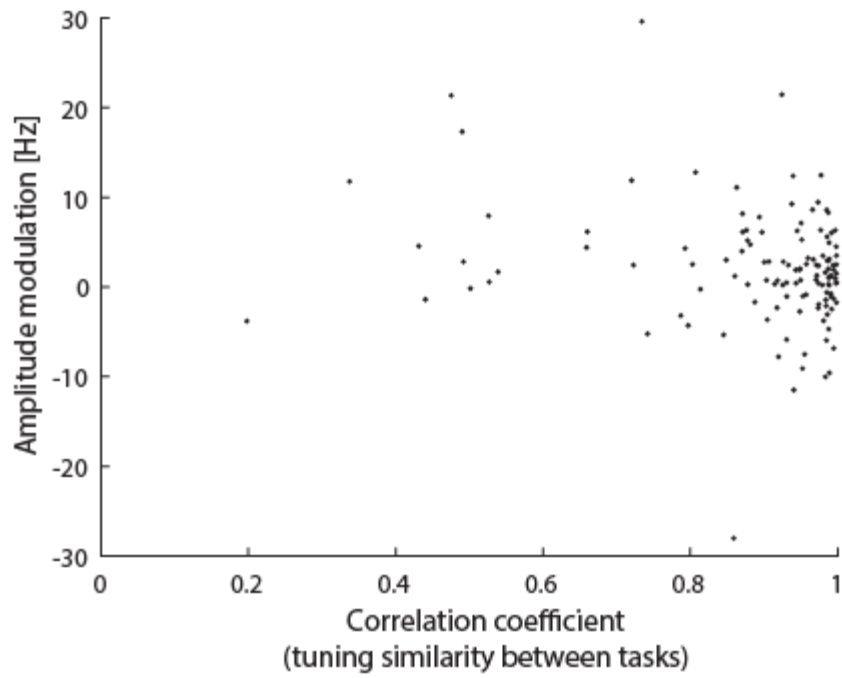
### Supplementary Figure S1.

Modulation of the population activity during the categorization and discrimination tasks for all visual stimuli (color 1–11). The conventions follow **Fig. 5a**.



**Supplementary Figure S2.**

The temporal evolution of population peak difference between fixation and discrimination tasks. The presentation conventions follow **Fig. 5c**.



### Supplementary Figure S3.

The response amplitude modulation (categorization – discriminaiton) is compared with the tuning curve correlation between the tasks (categorization vs. discrimination). Each dot indicates a single neuron. The correlation coefficient (the horizontal axis) and the amplitude modulation (the vertical axis) are the same as the ones used in **Figs. 2b** and **2c**, respectively.