# nature communications



**Article** 

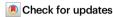
https://doi.org/10.1038/s41467-023-41743-3

# Sampling-based Bayesian inference in recurrent circuits of stochastic spiking neurons

Received: 31 January 2022

Accepted: 15 September 2023

Published online: 04 November 2023



Wen-Hao Zhang lacktriang 1,2,3,4,11, Si Wu $^{5,6,7,8}$ , Krešimir Josić lacktriang 9,10,12  $\boxtimes$  & Brent Doiron lacktriang 1,2,3,4,12  $\boxtimes$ 

Two facts about cortex are widely accepted: neuronal responses show large spiking variability with near Poisson statistics and cortical circuits feature abundant recurrent connections between neurons. How these spiking and circuit properties combine to support sensory representation and information processing is not well understood. We build a theoretical framework showing that these two ubiquitous features of cortex combine to produce optimal sampling-based Bayesian inference. Recurrent connections store an internal model of the external world, and Poissonian variability of spike responses drives flexible sampling from the posterior stimulus distributions obtained by combining feedforward and recurrent neuronal inputs. We illustrate how this framework for sampling-based inference can be used by cortex to represent latent multivariate stimuli organized either hierarchically or in parallel. A neural signature of such network sampling are internally generated differential correlations whose amplitude is determined by the prior stored in the circuit, which provides an experimentally testable prediction for our framework.

In an uncertain and changing world, it is imperative for the brain to reliably represent and interpret external stimuli. The cortex is essential for the representation of the sensory world, and it is believed that populations of neurons collectively code for richly structured sensory scenes<sup>1</sup>. However, two central characteristics of cortical circuits remain to be properly integrated into population coding frameworks. First, neuronal activity in sensory cortices is often noisy, showing significant variability of spiking responses evoked by the same stimulus<sup>2,3</sup>. In many traditional coding frameworks such spiking variability degrades the representation of stimuli by cortical activity<sup>4</sup>. Why cortical responses display large spiking variability while isolated cortical neurons can respond reliably remains far from clear. Second, the primary source of synaptic inputs to cortical neurons does not come from upstream centers which convey sensory signals, but rather from

recurrent pathways between cortical neurons<sup>5–7</sup>. While such recurrent connections are often organized about a stimulus feature axis<sup>8,9</sup>, it is not obvious whether or how their presence improves overall representation. We propose a biologically motivated inference coding scheme where these two ubiquitous cortical circuit features, variability in spike generation and recurrent connections, together support a probabilistic representation of stimuli in rich sensory scenes.

Numerous studies have framed sensory processing in the cortex in terms of Bayesian inference (e.g., refs. 10–16). Specifically, the 'Bayesian brain' hypothesis posits that sensory cortex infers and synthesizes a posterior distribution of the latent stimuli which describes the probability of possible stimuli that could have given rise to the sensory inputs. Performing Bayesian inference requires cortex to store an internal model that represents how sensory inputs and external

<sup>1</sup>Department of Neurobiology and Statistics, University of Chicago, Chicago, IL, USA. <sup>2</sup>Grossman Center for Quantitative Biology and Human Behavior, University of Chicago, Chicago, Chicago, IL, USA. <sup>3</sup>Department of Mathematics, University of Pittsburgh, Pittsburgh, PA, USA. <sup>4</sup>Center for the Neural Basis of Cognition, Pittsburgh, PA, USA. <sup>5</sup>School of Psychological and Cognitive Sciences, Peking University, Beijing 100871, China. <sup>6</sup>IDG/McGovern Institute for Brain Research, Peking University, Beijing 100871, China. <sup>7</sup>Peking-Tsinghua Center for Life Sciences, Peking University, Beijing 100871, China. <sup>8</sup>Center of Quantitative Biology, Peking University, Beijing 100871, China. <sup>9</sup>Department of Mathematics, University of Houston, Houston, TX, USA. <sup>10</sup>Department of Biology and Biochemistry, University of Houston, Houston, TX, USA. <sup>11</sup>Present address: Lyda Hill Department of Bioinformatics, UT Southwestern Medical Center, Dallas, TX, USA. <sup>12</sup>These authors contributed equally: Krešimir Josić, Brent Doiron. ⊠e-mail: kresimir.josic@gmail.com; bdoiron@uchicago.edu

stimuli are generated. Once a sensory input is received, cortical dynamics inverts this internal model in a process termed "analysis-by-synthesis"<sup>12</sup>, and represents the posterior distributively across neurons and/or across time<sup>15,16</sup>. In this study, we propose that recurrent connections in cortical circuits store the prior of latent stimuli to produce the posterior distribution when combined with evidence from sensory inputs. Moreover, we posit that Poisson spiking variability provides a source of fluctuations needed for generating random samples from the inferred posterior.

To test these hypotheses, we consider a recurrent circuit model where neurons receive stochastic feedforward inputs which carry information about the external world, and respond with Poissondistributed spiking activity. We find that such Poissonian spiking provides the variability that allows the network to generate samples from posterior stimulus distributions with differing uncertainties. We use this sampling framework to illustrate circuit-based Bayesian inference given two distinct generative models of stimuli in the external world: one organized hierarchically with a stimulus variable that depends on a latent stimulus parameter, and a second where a pair of latent stimuli are organized in parallel. In both cases, a recurrent circuit is able to generate samples from the joint posterior, and infer the values of the latent variables. We show through both analytic derivation and simulations that recurrent connections represent the correlation structure of these models, and the weight of these connections can be tuned to optimally capture the prior distribution of stimuli in the external world. The stronger the correlation between the latent variables, the stronger the recurrent connections need to be for the network to generate samples from the correct posterior distribution.

Finally, a neural signature of this circuit-based sampling mechanism is internally generated population noise correlations aligned with the stimulus response direction, often referred to as "differential correlations". In our framework, the amplitude of internally generated differential correlations is determined by the recurrent connection strength, which also determines the prior stored by the circuit. Since optimal inference requires a specific magnitude of recurrent connectivity, differential correlations resulting from such recurrent connectivity are a potential signature of optimal coding. This is in contrast to the deleterious impact of externally generated differential correlations. We thus predict that the correlation structure of the external world shapes recurrent wiring in neural circuits, and is reflected in the pattern of differential noise correlations. We use this logic to provide testable predictions from our framework for sampling-based Bayesian inference by recurrent, stochastic cortical circuits.

## **Results**

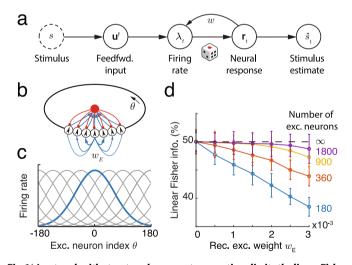
# Recurrent circuitry and spiking variability do not improve conventional neural codes

We start with the classic example of a sensory stimulus, s, encoded in neuronal population activity,  $\mathbf{r}$ , from which a stimulus estimate  $\hat{s}$  can be decoded (Fig. 1a, top)<sup>18</sup>. It is reasonable to expect that neuronal circuitry is adapted to accurately represent ethologically relevant stimuli. However, as we will show next, in simple coding schemes two ubiquitous features of cortical circuits – internal spiking variability and recurrent connectivity – are at best irrelevant for, and in many cases degrade, the accuracy of these representations.

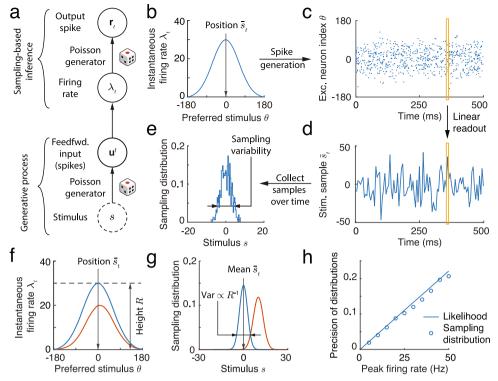
In population coding frameworks stimuli are encoded by a neuronal population with individual neurons tuned to a preferred stimulus value. The preferred values of all neurons cover the whole range of stimuli<sup>18-20</sup> (Fig. 1b, bottom); if *s* ranges over a periodic domain (such as the orientation of a bar in a visual scene, or the direction of an arm reach) then it is commonly assumed that the neurons' preferred stimuli are distributed on a ring (Fig. 1b, top). To generate neuronal responses from such a population we simulate a network of neurons whose spiking activity,  $\mathbf{r}_t$ , at time t is Poissonian with instantaneous firing rate  $\lambda_t$  (Eq. (11)). For simplicity we assume linear (or linearized) neuronal transfer

and synaptic interactions (Eqs. (10), (11)), so that the firing rate is a linear function of the feedforward and recurrent inputs. We couple excitatory (E) neurons with similar stimulus preferences more strongly<sup>8,9</sup> to one another, compared to neuron pairs with dissimilar tuning. In this way, the recurrent E connectivity has the same circular symmetry as the stimulus (Fig. 1b). In contrast, connections between inhibitory (I) neurons are unstructured, and inhibitory activity acts to stabilize network activity<sup>21</sup>. A stimulus, e.g., s = 0, results in elevated activity of E neurons with the corresponding preference (Fig. S1a). As expected, an increase in the strength of recurrent excitatory connections increases both the firing rates and the trial-to-trial pairwise covariability (i.e., noise correlations) in the responses<sup>2</sup> (Fig. S2a). This canonical network model has been widely used to explain cortical network dynamics and neural coding<sup>21-23</sup>. And our network model can produce neuronal responses that are qualitatively similar to experimental observations, including the variance of neuronal firing rate, the Fano factor, and the noise correlations (Fig. S2b-d).

We use linear Fisher Information (LFI) to quantify the impact of recurrent connectivity and internal spiking variability on the accuracy of the stimulus estimate,  $\hat{s}_t$ , from the activity vector  $\mathbf{r}_t$  (see details in Eq. S39 in Supplementary Information). The inverse of LFI provides a lower bound on the expected square of the difference between the true value, s, and the estimate,  $\hat{s}_t$ , made by a linear decoder LFI. In the limit of an infinite number of neurons available to the decoder LFI is unaffected by recurrent connectivity strength,  $w_E$  (Fig. 1d, dashed line). This is because the mean response of the network is linear in its inputs, and an (invertible) linear transformation can neither increase nor decrease LFI (see Eq. S38 in Supplementary Information). For networks with a finite number of neurons, the variability from spike generation is shared between neurons via recurrent interactions. Consequently, an increase in coupling strength,  $w_E$ , reduces LFI in finite networks (Fig. 1d, colored lines).



**Fig. 1**| A network with structured recurrent connections limits the linear Fisher Information (LFI) about external stimuli. a A schematic diagram showing how a stimulus, s, is encoded in neuronal response,  $\mathbf{r}_c$ . A stimulus estimate,  $\hat{s}_t$ , can be obtained from  $\mathbf{r}_t$ . **b** A recurrent ring model (top) where the connections between excitatory neurons are dependent on their distance along the ring. Blue arrows: excitatory synapses with line width denoting connection strength; red arrows: inhibitory synapses. **c** The population activity of excitatory neurons in the ring model,  $\mathbf{r}_t$ , dependent on a stimulus, s. The blue curve shows the population activity in response to s = 0, and gray curves the activities in response to stimuli with values at the peak locations of the curves. **d** For finite size networks (colored lines; ratio of excitatory to inhibitory neurons kept constant) LFI decreases as  $w_E$  increases. In the limit of infinite network size LFI does not depend on  $w_E$  (dashed line). Since neural responses are variable, LFI in the neuronal response converges to only half of the LFI in the feedforward input. Error bars denote one standard deviation (SD), which were estimated from N = 50 independent samples generated by using Bootstrap.



**Fig. 2** | **Spike generation with Poissonian variability can support sampling-based Bayesian inference. a** We use a feedforward network model (no recurrent connections) to demonstrate how spiking variability drives sampling. Neurons receive feedforward inputs,  $\mathbf{u}^t$ , modeled as independent Poisson spike trains, resulting in a Poissonian population response,  $\mathbf{r}_t$ , with means determined by the instantaneous firing rate vector,  $\boldsymbol{\lambda}_t$ . ( $\mathbf{b}$ - $\mathbf{e}$ ) Demonstration of sampling via stochastic spike generation. A population of neurons with Gaussian tuning and firing rates  $\boldsymbol{\lambda}_t$  ( $\mathbf{b}$ ) generates a realization of a population response,  $\mathbf{r}_t$  ( $\mathbf{c}$ ). A sample from the posterior distribution of the stimulus ( $\mathbf{d}$ , orange box) can be linearly read out from the population response ( $\mathbf{c}$ , orange box).  $\mathbf{e}$  The sampling distribution is obtained by

collecting stimulus samples over time. The profile of population firing rates ( $\mathbf{f}$ ) determines the sampling distribution ( $\mathbf{g}$ ). The position of the population firing rate,  $\bar{\mathbf{s}}_t$ , determines the mean of the sampling distribution, and the variance of the sampling distribution is inversely proportional to the peak firing rate, R. We show two population activity profiles, one in blue and the other in orange, to illustrate these points.  $\mathbf{h}$  In an E-I network, the precision of the sampling distribution (the inverse of sampling variability) read out from E neurons increases with the height of firing rate, and is consistent with the likelihood directly read out from the feedforward input.

In sum, recurrent connectivity and spiking variability do not improve, and often degrade, stimulus representation in the network (as measured by LFI). Since synaptic coupling is biologically expensive, a network that most accurately and cheaply represents a stimulus is then one with no recurrent connections (i.e.,  $w_E = 0$ ) and minimal spiking variability. Nevertheless, connectivity in mammalian cortex is highly recurrent solution of these extensive recurrent connections between cortical neurons in information representation, and why are their responses so noisy?

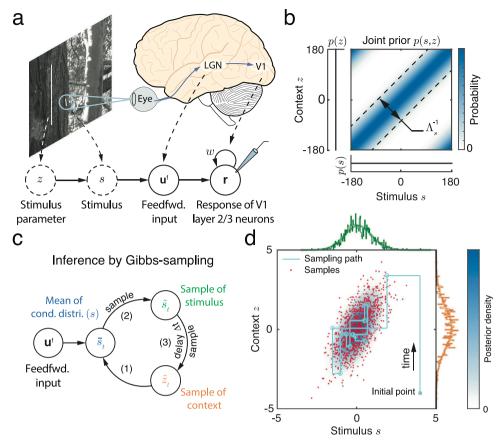
While classical population code theory often explains how to generate point estimates of a stimulus (Fig. 1a), numerous studies suggest that the brain performs Bayesian inference to synthesize and estimate the probability distribution of latent stimuli from sensory inputs (e.g., refs. 10–15,25,26). To compute this posterior a neural circuit needs to combine a stored representation of the prior distribution of the stimulus with the likelihood conveyed by feedforward inputs. We propose that recurrent connectivity can be used to represent the prior and spiking variability can generate samples from this posterior distribution. Before we present our full model we first show how sampling-based inference can be implemented in a population of spiking neurons.

# Internally generated Poisson spiking variability drives samplingbased Bayesian inference

Many studies suggest that neuronal response variability is a signature of sampling-based Bayesian inference in neural circuits (e.g.,

refs. 16,27-34). In these studies, the instantaneous population responses,  $\mathbf{r}_t$ , represent a sample of a latent stimulus, and the empirical distribution of stimulus samples collected over time is an approximation of the posterior distribution. Implementing sampling requires a network that generates variable output with stable statistics. It has been well documented that cortical spiking responses are often approximately Poissonian<sup>3,35</sup>. Theoretical studies suggest that such Poissonian variability can be internally generated in a network with dynamically balanced recurrent excitation and inhibition<sup>36,37</sup>. We thus assumed that our model neurons are Poissonian, and used the resulting fluctuations as the internal source of variability needed for sampling-based Bayesian inference. It remains to be shown if discrete Poissonian variability can be used to generate samples from stimuli with continuous probability distributions (e.g., orientation, moving direction) with the flexibility needed to represent different stimulus uncertainties. However, spike counts are discrete, and it is possible that errors that arise from representing continuous parameters by discrete random variable are characteristic of stimulus inference by animals that use sampling.

We address this question using a theory based on a simple model network composed of excitatory (E) Poissonian neurons (Eqs. (10), (11)), and subsequently support our findings by simulating a network containing both E and inhibitory (I) neurons (e.g., Fig. 1b). We start by showing that Poissonian spiking in a population of tuned neurons can drive sampling from a well–defined distribution. We assume that the instantaneous firing rates of a population of E neurons,  $\lambda_t$ , have a bell-shaped (Gaussian) profile (Fig. 2b), so that for the *j*th neuron



**Fig. 3** | **A hierarchical generative model and posterior inference via Gibbs sampling. a** An example of sensory feedforward input generation: The stimulus parameter, *z*, is the orientation of the tree trunk, and the stimulus, *s*, is the orientation of the bark texture located in the classical receptive field of a V1 hypercolumn. The recurrent circuit generates samples from the joint posterior over stimulus and stimulus parameter. Solid circles: observations and responses in the brain; dashed circles: latent variables in the external world. Nature image is adapted from Tkačik, G. et al. Natural images from the birthplace of the human eye. PLoS one 6, e20409 (2011). **b** The joint prior over the stimulus parameter, *z*, and stimulus,

s, is concentrated on the diagonal. The correlation between context and stimulus is determined by parameter  $\Lambda_s$ . (c) The posterior over stimulus parameter and stimulus can be approximated via Gibbs sampling (Eqs. (4a), (4c)) by iteratively generating samples of s and z from their respective conditional distributions. d The resulting approximations of the joint and marginal posterior over s and z. Light blue contour: the posterior distribution (Eq. (24)); Red dots: Samples obtained using Gibbs sampling. The green and orange projections are the marginal posterior distributions of s and z, respectively.

 $\lambda_{tj} = R \exp[\mathbf{h}_j(\bar{s}_t)] = R \exp[-(\bar{s}_t - \theta_j)^2/2a^2]$  (See Eq. (12) in Methods). Here  $\theta_j$  is the preferred stimulus of neuron j, a is the width of the tuning curve, and  $\bar{s}_t$  is the location of the peak of the firing rate profile,  $\lambda_t$ , in stimulus space (x-axis in Fig. 2b). Note that the value of  $\bar{s}_t$  is arbitrary here, but we will later relate it to the input to the population. Finally, the preferred stimuli of the E neurons,  $\{\theta_j\}_{j=1}^{N_E}$ , are uniformly distributed over the stimulus range (Fig. 1b). In each time interval the population activity is given by a vector of independent Poisson random variables,  $\mathbf{r}_t$ , with means determined by the instantaneous firing rate vector  $\lambda_t$  (Fig. 2b, c). At each time, t, this spiking activity produces a stimulus sample,  $\bar{s}_t$ , from the probability distribution determined by the instantaneous firing rates,  $\lambda_t$  (Fig. 2d, see Methods),

$$\tilde{s}_t \sim p(\tilde{s}|\boldsymbol{\lambda}_t) \propto \exp[\mathbf{h}(\tilde{s})^{\top}\boldsymbol{\lambda}_t] \propto \mathcal{N}(\tilde{s}|\bar{s}_t, \Lambda^{-1}).$$
 (1)

With the Gaussian firing rate profile we use here, the stimulus sample,  $\tilde{s}_t$ , can be read out as  $\tilde{s}_t = \sum_j \mathbf{r}_{tj} \theta_j / \sum_j \mathbf{r}_{tj}$  (Eq. (14) and Fig. 2d), which can be thought of as the location of the response,  $\mathbf{r}_t$ , in stimulus space (y-axis in Fig. 2c). The collection of stimulus samples across time ( $(\tilde{s}_t)$ ; Fig. 2e), determines the sampling distribution  $q(s) = T^{-1} \sum_t \delta(s - \tilde{s}_t)$  which approximates the distribution  $p(s|\lambda_t)$ , i.e.,  $p(s|\lambda_t) \approx q(s)^{16,38}$ . Here,  $\delta(\cdot)$  is the Dirac delta function and T is the number of samples. We assumed that instantaneous population firing rates are smooth to simplify the analysis, but this assumption is not essential. Sampling

driven by Poissonian variability will work as long as the temporally averaged population firing rate is smooth, even if the instantaneous population firing rate is noisy (see Eq. (17)).

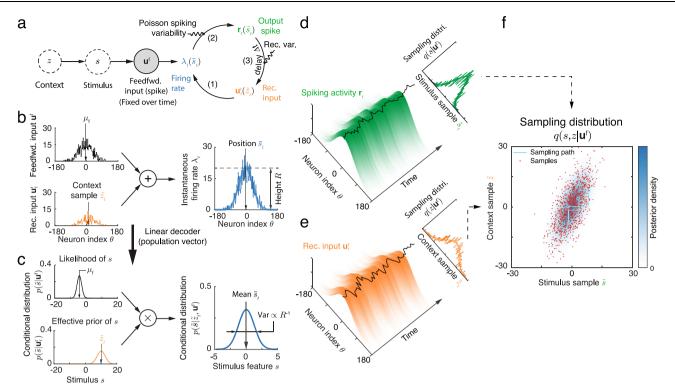
To use this mechanism to produce samples from the posterior distribution of a stimulus, we must define a generative model for the feedforward inputs evoked by a stimulus. We take the feedforward input to the neural population,  $\mathbf{u}^i$ , to be a vector of independent Poisson spike counts with Gaussian tuning over the stimulus, s. Following assumptions widely used in previous studies of probabilistic population codes (PPC)<sup>39,40</sup>, we assume that the mean input spike count to the jth excitatory neuron in the population is  $\langle \mathbf{u}^i_j(s) \rangle \propto \exp[\mathbf{h}_j(s)] = \exp[-(s-\theta_j)^2/2a^2]$ . A single realization of the input,  $\mathbf{u}^i$ , in a time interval encodes the whole likelihood function over the stimulus,  $p(\mathbf{u}^i|s)$ <sup>39</sup>. This likelihood is proportional to a Gaussian due to the Gaussian profile of feedforward input (Eq. (19)),

$$p(\mathbf{u}^{\mathsf{f}}|s) = \prod_{j=1}^{N_{\mathcal{E}}} \mathsf{Poisson} \left[ \langle \mathbf{u}_{j}^{\mathsf{f}}(s) \rangle \right],$$

$$\propto \exp \left[ \mathbf{h}(s)^{\mathsf{T}} \mathbf{u}^{\mathsf{f}} \right],$$

$$\propto \mathcal{N}(s|\mu_{\mathsf{f}}, \Lambda_{\mathsf{f}}^{-1}).$$
(2)

Here the likelihood mean,  $\mu_t$ , is determined by the location of  $\mathbf{u}^t$  in stimulus space, and the precision,  $\Lambda_t$ , is proportional to the spike count



**Fig. 4** | A recurrent circuit generates samples from the posterior defined by a hierarchical generative model. a Schematic of recurrent circuit dynamics, in which stimulus, s, and stimulus parameter, z, are encoded respectively in the population response,  $\mathbf{r}_t$ , and recurrent inputs,  $\mathbf{u}_t^r$ .  $\mathbf{b}$ ,  $\mathbf{c}$  When the feedforward inputs and recurrent inputs share the same tuning profile, summing the two inputs to define the instantaneous firing rate ( $\mathbf{b}$ ) is equivalent to multiplying the conditional distributions encoded by the two inputs to obtain the conditional distribution of the stimulus,  $p(s|\tilde{z}_t,\mathbf{u}^t)$ .  $\mathbf{c}$  The conditional distributions of the stimulus can be explicitly read out from corresponding population responses by a linear decoder ( $\mathbf{b}$ ).  $\mathbf{d}$ - $\mathbf{f}$ ) Reading out the joint sampling distribution from the recurrent circuit. The

projection of the spiking activity (Eq. (14)) and recurrent inputs (Eq. (29)) onto the stimulus subspace (black curves), can be read out linearly from the population activity and interpreted as a sample of stimulus and stimulus parameter respectively (Eqs. (4b), (4c)). Top right insets: the empirical marginal distributions of samples and marginal posteriors (smooth lines). (f) The joint value (red dots) of instantaneous samples of stimulus (black curve on the surface in (d)), and stimulus parameter (black curve on the surface in (e)) represent samples from the joint posterior of the stimulus and stimulus parameter. The true joint posterior is represented by the blue contour.

(or height) of  $\mathbf{u}^{\mathsf{f}}$  (Eq. (20)). Since a realization of the feedforward input encodes the whole likelihood function, we present a fixed  $\mathbf{u}^{\mathsf{f}}$  to the network over time (dropping the time index t), and describe how samples from the posterior  $p(s|\mathbf{u}^{\mathsf{f}})$  are generated by the network.

A simple example of inference via sampling is provided by a population of E neurons without recurrent connections and instantaneous firing rates equal to the feedforward input,  $\lambda_t = \mathbf{u}^t$  (Eq. (10)), and hence constant in time (Fig. 2a). In this feedforward network Poisson spike generation produces samples from the normalized likelihood, i.e.,  $\tilde{s}_t \sim p(\tilde{s}|\lambda_t) \propto p(\mathbf{u}^t|\tilde{s})$ , and consequently the network represents a uniform stimulus prior (i.e., p(s) is a constant).

To test our theory, we simulated the response of a network of tuned excitatory (E) and untuned inhibitory (I) neurons (Fig. 2a, c) to a fixed but randomly generated feedforward input (Eq. (18)). While the E neurons shared no recurrent connections, the E and I neurons were connected to maintain stable network activity. To confirm that the overall firing rate dictated the sampling variability (Eq. (1)), we increased the feedforward input rate, which reduced the width of the likelihood (Eq. (2)). As a result, the sampling precision (inverse of the sampling variance) increased and matched the precision of the likelihood (Fig. 2g, h), even as the normalized response variability (measured the by Fano factor) of single neurons remained unchanged.

While the above analysis introduces the key components of a sampling-based theory of inference, stimulus sampling using a feed-forward network is unnecessary: A single observation of the response  $\mathbf{r}$  in a deterministic feedforward network ( $\mathbf{r} = \mathbf{u}^t$  after removing spike generation in Eq. (11)) would also represent the whole likelihood<sup>39</sup>,

avoiding the costly process of collecting samples  $\tilde{s}_t$  across time. We next consider more interesting cases, and show that spiking variability in recurrent networks can drive sampling from more complex posterior distributions.

# Recurrent cortical circuit samples a hierarchical generative model

Recurrent networks can store a variety of generative model structures; to demonstrate the generality of our sampling framework we provide two example generative models which serve as building blocks for more complex models. We first consider a two-stage hierarchical model of feedforward inputs received by the cortical circuit (Fig. 3a). The first stage of our model consists of a stimulus, s, and a stimulus parameter, z, both of which are one dimensional for simplicity. The structure of the world is described by the joint distribution, p(s, z). Using the visual system as motivation, s, could be the orientation of the visual texture within a classical receptive field (local information) of a hypercolumn of V1 neurons, while stimulus parameter, z, may refer to the context orientation within a nonclassical receptive field of these cells (Fig. 3a). The likelihood of the stimulus based on a given parameter,  $p(s|z) = \mathcal{N}(s|z, \Lambda_s^{-1})$ , is Gaussian with precision  $\Lambda_s$ . For simplicity, we assume that the prior, p(z), is uniform, which implies that the marginal prior of s, is also uniform (Fig. 3b). This assumption is not essential for our main conclusions but does simplify the analysis. Importantly, the joint prior of stimulus and stimulus parameter, p(s, z), can have non-trivial structure with the density concentrated around the diagonal s=z (Fig. 3b). The precision  $\Lambda_s$  measures how strongly z and s are related, and thus determines how strongly their joint distribution is concentrated around the diagonal.

The second stage of the generative model describes how the feedforward input depends on the stimulus, s; this is identical to our prior treatment (See Eq. (2)). Combining these two stages provides a complete description of the generative model for the feedforward input received by neurons in the population,

$$p(\mathbf{u}^{\mathsf{f}}|s)p(s|z)p(z) \propto \prod_{j=1}^{N_{E}} \mathsf{Poisson}\Big(\mathbf{u}_{j}^{\mathsf{f}}|s\Big)p(s|z),$$
$$\propto \mathcal{N}(s|\mu_{\mathsf{f}},\Lambda_{\mathsf{f}}^{-1})\mathcal{N}(s|z,\Lambda_{\mathsf{s}}^{-1}). \tag{3}$$

Given this hierarchical model, we can show that the joint posterior over stimulus and stimulus parameters,  $p(s, z|\mathbf{u}')$  is a bivariate normal distribution (see Eq. (24)), and we next use it to evaluate the accuracy of the sampling distribution.

Gibbs sampling of the joint posterior of stimulus and stimulus parameter. One approach to approximate the joint distribution over stimulus and stimulus parameter is Gibbs sampling  $^{31,38,41,42}$  which starts with an initial guess for the value of the two latent variables, and proceeds by alternately generating samples of one variable from the distribution conditioned on the value of the second variable. More precisely, to approximate the joint posterior of s and z (Eq. (3)), Gibbs sampling proceeds by generating a sequence of samples,  $(\tilde{s}_t, \tilde{z}_t)$  indexed by time t, through recursive iteration of the following steps (Fig. 3c and Eq. (25)),

Compute : 
$$p(\tilde{s}|\tilde{z}_t, \mathbf{u}^f) \propto p(\mathbf{u}^f|\tilde{s})p(\tilde{s}|\tilde{z}_t) \equiv \mathcal{N}(\tilde{s}|\bar{s}_t, \Lambda^{-1}),$$
 (4a)

Sample : 
$$\tilde{s}_t \sim p(\tilde{s}|\tilde{z}_t, \mathbf{u}^t)$$
, (4b)

Sample: 
$$\tilde{z}_{t+\Lambda t} \sim p(\tilde{z}|\tilde{s}_t) = \mathcal{N}(\tilde{z}|\tilde{s}_t, \Lambda_s^{-1}).$$
 (4c)

Here  $\Delta t$  is the time increment between successive samples. The samples (red dots in Fig. 3d) are generated by alternately fixing the values of the two variables, so that sampling trajectories alternate between horizontal and vertical jumps (cyan lines in Fig. 3d). The empirical distribution of samples, i.e.,  $q(s,z|\mathbf{u}^f) = T^{-1} \sum_t \delta[(s,z)^\top - (\bar{s}_t,\bar{z}_t)^\top]$  with  $\top$  denoting vector transpose, approximates the joint posterior  $p(s,z|\mathbf{u}^f)$  (blue contour map in Fig. 3d, Eq. (24))<sup>38</sup>. To approximate  $p(s|\mathbf{u}^t)$ , the marginal posterior distribution of s, we can use only samples  $\bar{s}_t$  to obtain the approximating distribution  $q(s|\mathbf{u}^t)$  (compare the two green lines at the margin in Fig. 3d). The same is true for the marginal posterior over z.

Implementing Gibbs sampling of stimulus and stimulus parameter in a recurrently coupled cortical circuit. An implementation of Gibbs sampling in a recurrent E circuit can be intuitively understood by comparing the recurrent network dynamics (Fig. 4a) with the dynamics described by the Gibbs sampling algorithm (Fig. 3c). In the recurrent network a stimulus sample,  $\tilde{s}_t$ , is represented by the activity of E cells,  $\mathbf{r}_t$ , while a stimulus parameter sample,  $\tilde{z}_t$ , is represented by recurrent inputs,  $\mathbf{u}_t^r$ . To generate correct samples we require that the conditional distribution that is represented by the instantaneous firing rate,  $\lambda_t$  (Eq. (1)), matches the conditional distribution used in the Gibbs sampling algorithm (Eq. (4b)), so that  $p(\bar{s}|\tilde{z}_t,\mathbf{u}^t)=p(\bar{s}|\lambda_t)\propto \exp[\mathbf{h}(\bar{s})^T\lambda_t]$ . Equating the two distributions (see Eqs. (4a) and (10)) yields the relation,

$$\ln p(\tilde{s}|\tilde{z}_t, \mathbf{u}^f) = \ln p(\mathbf{u}^f|\tilde{s}) + \ln p(\tilde{s}|\tilde{z}_t),$$

$$\iff \mathbf{h}(\tilde{s})^\top \boldsymbol{\lambda}_t = \mathbf{h}(\tilde{s})^\top \mathbf{u}^f + \mathbf{h}(\tilde{s})^\top \mathbf{u}^t.$$
(5)

This equation holds when two constraints are satisfied: First, the firing rate vector,  $\lambda_t$ , needs to have a Gaussian profile peaked at  $\bar{s}_t$ , i.e., the mean of  $p(\tilde{s}|\tilde{z}_t,\mathbf{u}^f)$  (Eq. (4a)). Second, the peak firing rate, R, needs to be proportional to the precision of  $p(\tilde{s}|\tilde{z}_t, \mathbf{u}^f)$ , i.e.,  $R \propto \Lambda$  (see Fig. 2f, g). In a neural circuit one way for  $\lambda_t$  to satisfy these constraints is for feedforward inputs, u<sup>f</sup>, and recurrent inputs, u<sup>f</sup>, to both have Gaussian profiles with the same width,  $a_t$  as that of  $\lambda_t$  (by sharing the same  $\mathbf{h}(\tilde{s})$ , Egs. (5) and (12)). This is because the sum of two Gaussian-profile inputs with the same width, a, gives a firing rate,  $\lambda_t$ , with the same tuning, as long as the difference of the locations of two inputs is much smaller than the width, a. Our generative model (Eq. (3)) produces feedforward input, u<sup>f</sup>, with a Gaussian profile and encodes the likelihood function  $p(\mathbf{u}^f|\tilde{s})$ . The recurrent input,  $\mathbf{u}_r^r$ , then need to represent the conditional distribution  $p(\tilde{s}|\tilde{z}_t)$ . Hence, to satisfy Eq. (5) the recurrent input  $\mathbf{u}_t^r$  should have the same Gaussian profile as  $\mathbf{u}^t$  (Eq. (29)), with its location and magnitude determined by the mean and precision of  $p(\tilde{s}|\tilde{z}_t)$ , respectively.

If recurrent interactions are absent (setting  $\mathbf{u}_t^r = 0$ ), then network activity,  $\mathbf{r}_b$  generates samples from the normalized likelihood,  $p(\mathbf{u}^f | \tilde{s})$ , as we showed previously when describing feedforward networks (Fig. 2). When neurons only receive recurrent inputs (setting  $\mathbf{u}^f = 0$ ), the network generates samples from the conditional distribution  $p(\tilde{s}|\tilde{z}_t)$ . Driven by a sum of recurrent and feedforward inputs, the network generates samples from a distribution given by the product of the conditional distributions encoded by both inputs respectively (Fig. 4b, c).

The recurrent weights must be adjusted so that the recurrent input has the appropriate magnitude and width to encode the likelihood p(s|z). To simplify the exposition we first assume that E neurons are only self-connected, so that the width of recurrent input trivially matches that of the feedforward input (otherwise recurrence will broaden the profile of the firing rate activity  $\lambda_t$  over the network). To constrain the magnitude of the recurrent weights we require that the sum of the recurrent inputs satisfies  $\sum_j \mathbf{u}_{tj}^r \propto \Lambda_s$ . Since  $\mathbf{u}_j^r = w_E \mathbf{r}_j$  and the width of  $\mathbf{u}_j^r$  and  $\mathbf{r}_j$  are equal, the magnitude of the recurrent weights that result in samples from the correct posterior must satisfy:

$$\boldsymbol{w}_{E}^{*} = \frac{\langle \mathbf{u}_{j}^{r} \rangle}{\langle \mathbf{r}_{i} \rangle} = \frac{\langle \sum_{j} \mathbf{u}_{j}^{r} \rangle}{\langle \sum_{i} \mathbf{r}_{i} \rangle} = \frac{\Lambda_{s}}{\Lambda_{f} + \Lambda_{s}},$$
 (6)

where  $\Lambda_s$  and  $\Lambda_f$  are the precision of likelihood p(s|z) and  $p(\mathbf{u}^f|s)$ respectively (Eq. (3)). The optimal recurrent weight,  $w_F^*$ , thus encodes the correlation between the stimulus s and the stimulus parameter z. An increase in correlation between s and z, resulting in a narrower diagonal band in p(s, z) (Fig. 3b), requires an increase in the recurrent weight  $w_F^*$  for optimal sampling. When the underlying parameter and stimulus are uncorrelated so that  $\Lambda_s = 0$ , the hierarchical generative model (Fig. 3a) is equivalent to the generative model without stimulus parameter (Fig. 2a) and recurrent interactions are not needed for sampling (i.e.,  $w_F^* = 0$ ). Moreover, the optimal recurrent weight also depends on the likelihood precision  $\Lambda_f$  that is determined by the input spike count. Hence, the optimal weight needs to be adjusted depending on feedforward inputs so that samples from the correct posterior are generated (see Discussion of how this feature impacts the network sampling). Overall, our framework (Eq. (6)) thus predicts that optimal Bayesian inference is achieved with recurrent synaptic weights which depend on the correlative structure of the external world. We numerically test this prediction in the next section.

# A stochastic E-I spiking network jointly samples stimulus and stimulus parameter

To confirm the predictions of this analysis, we simulated a full recurrent network consisting of both E and I neurons with Poisson spike train statistics (see details in Eqs. (47)–(50)). The E neurons were synaptically connected to each other (Eq. (49), see Fig. 1a), in contrast

to the simple network of self-connected E neurons we described above. While recurrent E to E coupling broadens the tuning of excitatory recurrent input, lateral inhibition can sharpen Gaussian firing rate profiles so that it matches that of the feedforward inputs (as required by Eq. (5)).

The activity of the recurrent network in response to a fixed but randomly generated feedforward input (Eq. (3)) can be decoded to produce samples from the bivarite posterior distribution of the stimulus and stimulus parameter. As above, samples from the conditional stimulus distribution are represented by the activity of E neurons (Eq. (14)), while samples from the conditional stimulus parameter distribution are represented by recurrent inputs received by E neurons (Eq. (29); black curves overlaid on the top of population responses in Fig. 4d, e, respectively). To update recurrent inputs we only used neuronal activity at the previous time step. Thus, the activities of E neurons and their recurrent inputs were updated in alternation, consistent with Gibbs sampling. The trajectory obtained by plotting the stimulus sample read out from the network activity on one axis, and plotting the stimulus parameter sample read out from recurrent E inputs on another axis then exhibits the characteristics of Gibbs sampling (Fig. 4f, cyan line). The resulting sampling distribution provides a good approximation to the joint posterior of stimulus and context (compare red dots and blue contour in Fig. 4f). Inhibitory neurons again did not respond selectively to either the stimulus or the stimulus

For the network to generate samples from the joint posterior, the recurrent connectivity should depend on the correlation between the stimulus and the stimulus parameter (Eq. (6)). To verify this prediction, we fixed the generative model (Eq. (3)) and changed only the recurrent weights in the network. For simplicity, we only varied the peak E weight,  $w_E$  (Eq. (49)), and maintained network stability by fixing the ratio between E and I synaptic weights. While increasing  $w_E$  did not change the sampling mean, it did increase the variance of the stimulus parameter sampling distribution, and increased the correlation between stimulus and stimulus parameter samples (Fig. 5a).

We use Kullback–Leibler (KL) divergence to measure the distance between the sampling distribution,  $q(s,z|\mathbf{u}^t)$ , and the true posterior,  $p(s,z|\mathbf{u}^t)$  (Eq. (24)). The KL divergence quantifies the loss of mutual information, measured in bits, between the latent variables (s and z) and the feedforward inputs,  $\mathbf{u}^t$ , when the true posterior, p, is approximated by the distribution, q (Eq. (42))<sup>38</sup>. The mutual information loss in the network is minimized at a unique value of the recurrent weight,  $w_E^t$ , at which the sampling distribution, q, best matches the posterior, p (Fig. 5b, black circle). To confirm that this optimal recurrent weight,

 $w_F^*$ , increases with the correlation in the prior (precision  $\Lambda_s$ , Eq. (6)), we numerically obtained the recurrent weight that minimizes the mutual information loss for each value of  $\Lambda_s$  in the generative model. These results confirmed the predictions of our theory (Eq. (6), Fig. 5c): When  $\Lambda_s = 0$ , i.e., when stimulus parameter and stimulus are uncorrelated, a network with no interactions performs best ( $w_E^* = 0$ ), while for small  $\Lambda_s$  (relative to  $\Lambda_t$ ) the optimal weight  $w_E^*$  is positive and increases with  $\Lambda_s$ . In total, we have described a potential mechanism for a recurrent network of spiking neurons to perform sampling-based Bayesian inference.

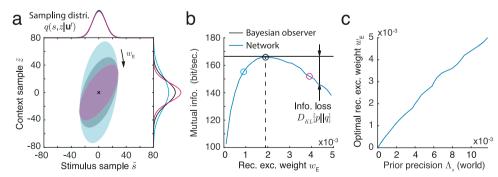
# Generating samples from multi-dimensional posteriors with coupled neural circuits

To demonstrate the generality of the proposed neural code we next consider a world described by a broad, rather than deep (hierarchical) generative model. Information about each of two latent stimuli,  $\mathbf{s} = (s_1, s_2)$ , is relayed by corresponding feedforward inputs received by a neural circuit (Fig. 6a). We assume the prior is a bivariate Gaussian distribution (Fig. 6b), i.e.,  $p(\mathbf{s}) \propto \exp[-\Lambda_s(s_1-s_2)^2/2] \equiv \mathcal{N}(s_1-s_2,\Lambda_s^{-1})$ , so that  $\Lambda_s$  ( $\Lambda_s \geq 0$ ) characterizes the correlation between  $s_1$  and  $s_2$ . Furthermore, each stimulus,  $s_m$ , independently generates feedforward spiking inputs,  $\mathbf{u}_m^f$ , each of which is received by a separate network and produces responses  $\mathbf{r}_m$  for m=1,2 (Fig. 6a). Thus, the full generative model of the input has the form,

$$p(\mathbf{u}^{\mathsf{f}}|\mathbf{s})p(\mathbf{s}) = \left[\prod_{m=1}^{2} p(\mathbf{u}_{m}^{\mathsf{f}}|s_{m})\right] p(s_{1},s_{2}),$$

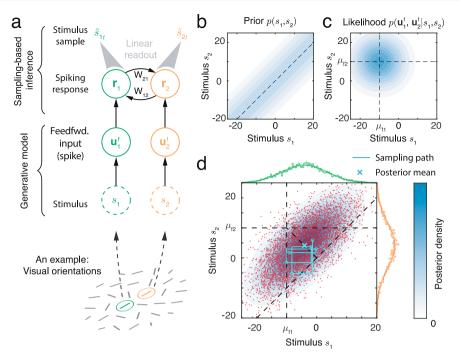
$$\propto \left[\prod_{m=1}^{2} \mathcal{N}(s_{m}|\mu_{\mathsf{f}m},\Lambda_{\mathsf{f}m}^{-1})\right] \mathcal{N}(s_{1} - s_{2},\Lambda_{\mathsf{s}}^{-1}).$$
(7)

The likelihood  $p(\mathbf{u}_m^f|s_m)$  is the same as that given previously (Eq. (2)), where the feedforward inputs,  $\mathbf{u}_m^f$ , are again described by conditionally independent Poisson spike counts with Gaussian tuning over stimulus  $s_m$ . As a concrete example, the two stimuli,  $s_m$ , could represent orientations of local edges falling in the central receptive fields of a V1 hypercolumn (Fig. 6a, bottom), with each V1 hypercolumn modeled by a network producing the response  $\mathbf{r}_m$  (Fig. 6a, top). Then  $\Lambda_s$  characterizes a priori tendency of the stimuli to share similar orientations, and determines how likely two local edges are to be part of a global line, as in the case of contour integration 43,44. However, the generative model defined by Eq. (7) is quite general and has been also used to explain multisensory cue integration and sensorimotor learning 3.



**Fig. 5** | **The joint sampling distribution of stimulus and stimulus parameter changes with the recurrent weight in the network. a** The sampling distribution for different recurrent excitatory weights,  $w_E$ . The ratio of excitatory and inhibitory weights was fixed. Ellipses capture three standard deviations from the mean of the joint sampling distribution. Different colors correspond to the three values of  $w_E$ , denoted by different symbols in **b. b** The mutual information between the latent variables, s and s, and the feedforward inputs for an ideal Bayesian observer (black

horizontal line) and for the sampling distribution generated by the network model (blue curve). The difference between the two lines is the KL divergence between the posterior,  $p(s, z|\mathbf{u}^t)$ , and the sampling distribution,  $q(s, z|\mathbf{u}^t)$ . KL divergence is minimized when the weight in the recurrent network is set to a value,  $w_E^*$ , at which the sampling distribution, q, best matches the true posteriori, p (black circle).  $\mathbf{c}$  This optimal weight,  $w_E^*$ , increases with prior precision,  $\Lambda_S$ .



**Fig. 6** | **Distributed sampling from a multivariate posterior distributions using coupled networks.** a Network m (m = 1, 2) receives a feedforward input evoked by a stimulus,  $s_m$ . The coupling between the two networks represents the stimulus prior. A linear readout from each network, m, can be interpreted as a sample from the posterior of the stimulus,  $s_m$ . Examples of a prior (**b**) and likelihood (**c**). The prior distribution is concentrated around the diagonal line (dashed line), indicating the two stimuli are more likely to be colinear. In (**c**),  $\mu_{\Pi}$  = -10 and  $\mu_{I2}$  = 10 are the

means of the likelihoods of  $s_1$  and  $s_2$ , respectively. **d** The joint posterior of stimuli and the corresponding approximate sampling distribution generated by the coupled networks. A sample from the joint posterior can be read out individually from the activity of the corresponding network (shown in **a**). Light blue contour: the posterior distribution (Eq. (34)); Red dots: stimulus samples generated by the network.

The posterior is a bivariate Gaussian distribution (Fig. 6d, Eq. (34)) whose mean is shifted from the likelihood mean (Fig. 6c) towards to the diagonal line, because of the correlations between the stimuli in the prior (Fig. 6b). We can again use Gibbs sampling to approximate the posterior  $p(\mathbf{s}|\mathbf{u}^t)$  using the following steps,

Compute : 
$$p(\tilde{s}_1|\mathbf{u}_1^t, \tilde{s}_{2t-\Lambda t}) \propto p(\mathbf{u}_1^t|\tilde{s}_1)p(\tilde{s}_{2t-\Lambda t}|\tilde{s}_1)$$
, (8a)

Sample: 
$$\tilde{s}_{1t} \sim p(\tilde{s}_1 | \mathbf{u}_1^f, \tilde{s}_{2,t-\Delta t})$$
, (8b)

where  $\tilde{s}_{1t}$  and  $\tilde{s}_{2t}$  are instantaneous samples at time t of stimuli  $s_1$  and  $s_2$ , respectively. We only give the steps needed to produce samples from the conditional distribution of  $s_1$ , as samples from the conditional distribution of  $s_2$  can be obtained using the same steps after exchanging indices.

These sampling steps can be implemented distributively in a coupled neural circuit using a mechanism similar to that we described in the case of a hierarchical generative model. The activity of each network,  $\mathbf{r}_{m}$ , individually represents samples from the (marginal) posterior of  $s_m$  (Fig. 6a, top). The joint posterior is then approximated as the collection of samples represented by the activity pairs  $(\mathbf{r}_1, \mathbf{r}_2)$ . Taking network m=1 as an example, spike response  $\mathbf{r}_{1t}$  produces a stimulus sample  $\tilde{s}_{1t}$  as long as the instantaneous firing rate  $\lambda_{1t}$  represents the conditional distribution  $p(\tilde{s}_1|\mathbf{u}_1^f, \tilde{s}_{2,t-\Delta t})$  (Eq. (8a)). Since the feedforward input,  $\mathbf{u}_1^f$ , represents the likelihood  $p(\mathbf{u}_1^f | \tilde{s}_1)$ , to obtain the appropriate firing rates,  $\lambda_{1t}$ , the recurrent input from network 2 to network 1,  $\mathbf{u}_{12,t}^{r}$ , must encode the correct conditional distribution,  $p(\tilde{s}_{2,t-\Delta t}|\tilde{s}_1)$ . As in the case of the mechanism we proposed to implement sampling as described by Eq. (5),  $\mathbf{u}_{12,t}^{r}$  needs to have the same Gaussian profile as the firing rate  $\lambda_{1t}$ , the position of  $\mathbf{u}_{12,t}^{r}$  on the stimulus space should match the mean of  $p(\tilde{s}_{2,t-\Delta t}|\tilde{s}_1)$ , i.e.,

 $\tilde{s}_{2,t-\Delta t} = \sum_j \mathbf{u}_{12,tj}^r \theta_j / \sum_j \mathbf{u}_{12,tj}^r$ , and the magnitude of  $\mathbf{u}_{12,t}^r$  must be proportional to the prior correlation,  $\Lambda_s \propto \sum_j \mathbf{u}_{12,tj}^r$  (Eq. (39)). Hence, each network can sum the feedforward input and the recurrent input from its counterpart to obtain an update to the instantaneous conditional distribution given by Eq. (8a), and generate independent Poisson spikes to produce a sample from the instantaneous conditional distribution (Eq. (8b)). Notably, the sample of each stimulus can be locally read out from corresponding network (Eq. (41), Fig. 6a), even if the activities of two networks are correlated.

Since the recurrent input strength represents the stimulus correlation in the prior determined by precision  $\Lambda_s$ , the coupling between the two networks needs to be tuned to generate the appropriate recurrent input. Indeed, in a network with only E neurons, and connections only between neurons with the same preferred stimulus value but in different networks, the optimal homogeneous connection strength is  $w_{mn}^* = \langle \mathbf{u}_{mn,j}^r \rangle / \langle \mathbf{r}_{n,j} \rangle = \Lambda_s / (\Lambda_{fn} + \Lambda_s)$  (Eq. (40)). This mirrors the result obtained with the hierarchical model presented earlier in Eq. (6).

**Coupled E-I spiking networks sample bivariate dimensional posteriors.** To test the feasibility of the proposed mechanisms for generating samples from a bivariate posterior, we simulated a pair of bidirectionally coupled circuits consisting of E and I neurons (Fig. 7a). This neural circuit model can be extended to generate samples from higher dimensional posterior distribution (see Discussion). Each circuit receives feedforward input generated by one of the two stimuli. On every time step the sample of each stimulus,  $\tilde{s}_{mt}$ , can be individually and linearly read out from the response of corresponding network,  $\mathbf{r}_{mt}$  (Eq. (41)). Jointly, the two stimulus samples, one each from both networks,  $\tilde{\mathbf{s}}_t = (\tilde{\mathbf{s}}_{1t}, \tilde{\mathbf{s}}_{2t})^{\mathsf{T}}$ , provide a sample from the joint posterior of the two latent stimuli (Fig. 7b). We assumed that the synaptic connections between the networks,  $w_{mn}$   $(m, n=1, 2; m \neq n)$ , are excitatory, but

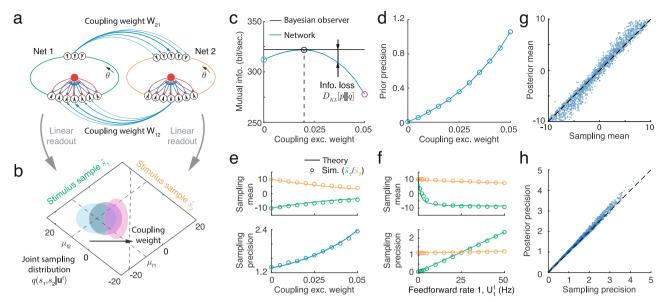


Fig. 7 | The statistics of the multivariate sampling distribution of stimuli generated by coupled E-I circuits. a Each of the two circuits individually generate a sample of a corresponding stimulus which can be read out linearly from that circuit's activity. Combining the readouts from the two networks yields the joint sampling distribution. The ring color indicates the stimulus sample the circuit generates: green and orange represent the stimulus  $s_1$  and  $s_2$ , respectively. Blue arrows: E synapses with width denoting connection strength; red arrows: I synapses. b The sampling distribution shifts from the likelihood mean to the diagonal line as the coupling between the networks increases. Ellipses capture one standard deviation from the mean of the sampling distribution. Different colors correspond to the three different coupling weights between the circuits shown in (c). c The mutual information between latent variables and the feedforward inputs for the ideal Bayesian observer (black) and the sampling distributions generated by the

network with different coupling weights between the two circuits. **d** The optimal coupling weight that minimizes information loss also increases with prior precision (which is inversely proportional to the width of the band in Fig. 6b). **e** The mean and precision of the sampling distribution over the two stimuli change with the coupling weight between the circuits when the feedforward input is fixed. **f** The mean and precision of the sampling distribution over the two stimuli change with the firing rate of feedforward input to network 1, with other network parameters fixed. Comparison of the mean (**g**) and precision (**h**) of the sampling distributions with the posteriors under different combinations of feedforward inputs and coupling weights. Different dots are obtained from the sampling distributions obtained under different combinations of input direction and strength, and coupling weight between networks.

target both E and I neurons, while inhibitory connections are local to each network. We also adjusted network parameters so that the profiles of the inputs across networks (e.g., the inputs from network 2 to 1) have the same tuning profile as the feedforward inputs (see Methods). Since we assumed uniform marginal priors (see Eq. (32)), recurrent connections between E neurons within the a circuit were absent, while E and I neurons within a circuit were recurrently connected to ensure network stability. For simplicity, we chose parameters so that the two circuits were symmetric, but the strength of the feedforward inputs to each could differ.

We asked whether the activity of the two coupled circuits can generate samples from bivariate posteriors, and how the sampling distribution depends on the coupling,  $w_{mn}$ , between the two circuits. An increase in synaptic coupling between the two networks caused the sampling distribution to shift from the likelihood mean towards the diagonal (Fig. 7b), resulting in stimulus samples,  $\tilde{s}_{1t}$  and  $\tilde{s}_{2t}$  that were more similar. This is consistent with an increase in stimulus correlation in the multivariate prior,  $\Lambda_s$  (Eq. (7)). To confirm our prediction that the optimal coupling strength between the two networks,  $w_{mn}^*$ , increases with the stimulus correlation in the prior,  $\Lambda_s$ , we numerically obtained the coupling weight that minimizes the loss of mutual information between latent stimuli and feedforward inputs (Fig. 7c). The optimal synaptic weight between the circuits increased with stimulus correlation in the prior. At the optimal weight,  $w_{mn}^*$ , the sampling distribution was close to the true posterior, showing that a properly tuned circuit can generate samples from the correct distribution (Fig. 7d).

We next asked how the sampling distribution in the network depends on network and feedforward input parameters. As the coupling between the two circuits increased, the sample means of both stimuli converge (Fig. 7e, top) and the sampling precision of both stimuli increased as well (Fig. 7e, bottom), in agreement with a more correlated stimulus prior. We also tested whether a network with fixed parameters can generate samples from a family of posteriors with different uncertainties. To do so, we changed the uncertainty of the likelihood of  $s_1$  by changing the firing rate in the feedforward input  $\mathbf{u}_1^f$ received by network 1. We observed that with a narrower likelihood of  $s_1$ , the sample means of both stimuli shifted towards the mean of likelihood of  $s_1$  (-10°), and sampling precision increased, consistent with a change in the posterior distribution (Fig. 7f). Lastly, to demonstrate the robustness of this network implementation of samplingbased inference we compare the sampling distributions to the true posteriors under different combinations of input and network parameters (Fig. 7g, h), in each case setting the recurrent coupling to the optimal value,  $w_{mn}^*$ , obtained numerically. Across different parameter values, we observe excellent agreement in both the mean (Fig. 7g) and precision (Fig. 7h) of the two densities. In sum, our recurrent network of spiking neuron models can be extended to support sampling-based Bayesian inference with multi-dimensional stimuli.

# A signature of stimulus sampling: internally generated differential noise correlations

A central prediction of our circuit framework for sampling-based Bayesian inference is that an increase in the correlation between stimuli in the sensory world should result in stronger synapses between neurons whose activities represent these stimuli (see Eq. (6)). This is a difficult prediction to test since measuring synaptic connectivity along a functional axis is already challenging<sup>45</sup>, let alone measuring a change in synaptic strength owing to a change in stimulus statistics. Here, we outline a testable prediction of our theory by identifying a measurable,

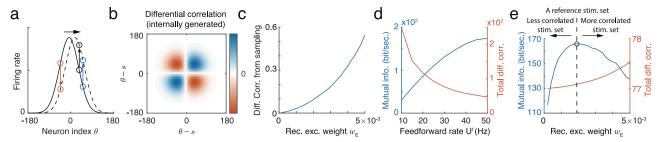


Fig. 8 | Stimulus sampling by a network is reflected in the internally generated differential correlations, whose impact differs from differential correlations inherited from feedforward inputs. a Stimulus sampling via spike generation causes the population firing rate to fluctuate along the stimulus subspace (x-axis). b The pattern of internally generated differential correlation in a network implementing sampling composed of neurons with Gaussian tuning. c Internally generated differential correlations in such a network increase with recurrent weight,  $w_E$ . d The rate in feedforward input decreases the externally generated correlations,

and increases the mutual information between the feedforward inputs and latent stimulus. **e** Recurrent network weights increase internally generated differential correlations. Mutual information between stimulus and feedforward inputs changes non-monotonically with recurrent weight. The direction of arrows indicates the predicted direction of change of the recurrent weights after an animal is retrained using a new stimulus set with different correlations compared to the reference stimulus set.

population-level signature of changes in functionally related recurrent synaptic strengths.

In response to a fixed feedforward input the responses of a recurrent circuit implementing stimulus sampling will fluctuate. The alignment of the recurrent circuitry and neuronal stimulus tuning causes a portion of these activity fluctuations to align with the subspace in which stimuli are coded. As an example, consider the sampling implemented by a single recurrent network (Fig. 4a), and suppose the population response fluctuates around its mean position (0° in the example of Fig. 8a), ignoring fluctuations along other directions in neuronal response space. The activity of neuron pairs with stimulus preference both above or below the mean position are positively correlated (the black and blue neurons in Fig. 8a), while the activity of neuron pairs with preferences straddling the mean are negatively correlated (the black and red neurons in Fig. 8a). Such stimulus sampling generates a covariance component which is proportional to the outer product of the derivative of neuronal tuning (Fig. 8b), i.e.,  $\mathbf{f}_{s}'\mathbf{f}_{s}'^{\mathsf{T}}$ , where  $\mathbf{f}'_s$  denotes the derivative of tuning  $\mathbf{f}(s) = \langle \boldsymbol{\lambda}_t \rangle$  (mean firing rate) over stimulus s. Such noise correlations have been referred to as differential correlations4,17, and are generally viewed as deleterious to stimulus coding. Stochastic sampling in coupled networks (Fig. 6a) produces similar differential noise correlations (see Supplementary Information).

In our network implementation of sampling, the amplitude of internally generated differential correlations is not arbitrary, but is determined by the recurrent connection strength,  $w_E^*$ . Here, the differential covariance matrix of population responses has the form (see Eq. (44))

$$\Sigma_{DC} = V(\bar{s}|\mathbf{u}^{f})\mathbf{f}_{s}'\mathbf{f}_{s}^{T},$$
where  $V(\bar{s}|\mathbf{u}^{f}) = \frac{\Lambda_{s}}{\Lambda_{f}(\Lambda_{f} + \Lambda_{s})} = a^{2}n_{f}^{-1}w_{E}^{*},$ 
(9)

where  $V(\bar{s}|\mathbf{u}^t)$  is the variance of  $\bar{s}_t$  in equilibrium over time, and  $\bar{s}_t$  is the mean of the instantaneous conditional distribution (Eq. (4a)) represented by the position of instantaneous firing rate  $\lambda_t$  (Fig. 2b). Importantly, the amplitude of differential correlations increases with the recurrent weight,  $w_E^*$ , which is set by the prior precision  $\Lambda_s$  (Eq. (6); Fig. 8c). Thus, in our framework internally generated differential correlations are a by-product of inference by sampling from posterior distributions of stimuli in a structured world.

**Distinguishing external and internal differential correlations.** The previous analysis of internally generated differential correlations in a circuit implementing sampling-based inference is based on the assumption of a fixed feedforward input (Eq. (9)). However, in typical

neurophysiology experiments an external stimulus, s, is fixed, while the feedforward input,  $\mathbf{u}^t$ , fluctuates due to variability in sensory acquisition and transmission noise (Eqs. (3) and (7)). Hence, differential correlations of neuronal population responses are a combination of correlations inherited from feedforward input<sup>46</sup>, and correlations generated by recurrent network interactions that align with the population stimulus tuning<sup>24</sup>. When the feedforward input is described by a hierarchical generative model (Eq. (2)), the total magnitude of differential correlations in the evoked response is  $a^2n_{\rm f}^{-1}w_{\rm f}\mathbf{f}_s^{\prime}\mathbf{f}_s^{\prime}\mathbf{f}_s^{\prime}\mathbf{f}_s^{\prime}\mathbf{f}_s^{\prime}$  (see Eq. (46)), where the second term reflects differential correlations inherited from the feedforward input (compare with Eq. (9)). Although the two sources of differential correlations are intertwined in the neuronal response, they impact the information content differently thus offering a potential way to distinguish between them in neural data.

Externally generated differential correlations decrease with feedforward input rate, which could be modulated by visual stimulus strength such as contrast (Fig. 8d, red curve). As a consequence, the mutual information (the information between feedforward inputs uf and the latent variables, i.e., s and z, sampled by recurrent network in Fig. 4a, Eq. (42)) increases with feedforward input intensity (Fig. 8d, blue curve). We, therefore, have a monotonic, decreasing relationship between externally generated differential correlations and mutual information. This is expected since such inherited correlations always impair information processing, as observed previously<sup>4,17</sup>. In contrast, an increase in recurrent weights,  $w_E$ , increases internally generated differential correlations, but results in a non-monotonic change in mutual information (Fig. 8b). Hence there is a non-monotonic relation between internally generated differential correlations and the mutual information between stimulus and feedforward inputs. In sum, the impact of external and internal differential correlations on stimulus coding can be distinguished by their respective monotonic and nonmonotonic relation with the mutual information between stimulus and response.

#### **Discussion**

We have presented a framework in which neuronal response variability and recurrent synaptic connections, two ubiquitous features of cortex, are jointly used to implement sampling-based Bayesian inference in neuronal circuit models. Combining mathematical analysis and network simulations, we established that stereotypical Poisson variability of discrete spike counts can drive flexible sampling from a family of continuous distributions. The sampling statistics are determined by the structure of recurrent coupling, which stores information about the stimulus prior, and feedforward inputs conveying the stimulus likelihood. Sampling-based inference is implemented in two steps: the

instantaneous firing rate, determined by the sum of feedforward and recurrent inputs, represents the instantaneous conditional distribution of latent stimulus, while Poissonian variability in spike generation is used to generate a random stimulus sample from this conditional distribution. We have shown how sampling can be implemented using biologically feasible mechanisms for three different generative models of increasing complexity. The simplest model includes one latent stimulus, while the more complex models include multiple latent stimuli organized hierarchically or in parallel. These three generative models form the basic building blocks of more complex models. Thus our ideas can be extended to a wide range of perceptual and cognitive processes<sup>47</sup>.

The neural code we described shares some features with codes described in previous studies, including parametric representations in probabilistic population codes (PPCs)<sup>15,39,40</sup>, and sampling-based codes (SBCs) <sup>16,27-32</sup>. In our framework, the conditional distributions of latent variables is represented by instantaneous firing rates which linearly encode the logarithms of these conditional distributions, and have a mathematical form that is similar to that used in past studies describing PPCs (e.g., Eq. (5)). Further, the posterior is represented by stimulus samples generated through a random process, a feature of all SBCs. Despite these similarities, there are fundamental differences between the neural code we described and previously proposed PPCs and SBCs.

PPCs are generally implemented in networks with no internally generated variability, with stochasticity inherited from the stimulus. In contrast, our proposed network is doubly stochastic: The Poisson variability in the feedforward input allows a single realization of the feedforward input to represent the whole stimulus likelihood<sup>39</sup>, while internally generated Poisson variability drives stimulus sampling. Further, in PPCs the posterior is represented parametrically by a oneshot neuronal response, while in our proposed network the joint posterior is approximated by a sequence of samples, each obtained as a linear readout from the instantaneous neuronal responses. Although it takes time to collect sufficiently many samples to approximate the posterior well, an advantage of sampling codes compared to PPCs is that inference with multivariate posteriors can be implemented using linearly coupled subnetworks (Fig. 6), with the number of subnetworks determined by the dimension of the latent stimulus features. In contrast, to represent an M-dimensional multivariate posterior using PPCs requires  $N^{M}$  neurons in a linear network (N is the number of neurons in representing each dimension) so that the number of neurons increases exponentially with the latent stimulus dimension,  $M^{16}$ . Alternatively, coupled networks with NM neurons can be used, but require complex, nonlinear coupling between these networks<sup>48,49</sup>.

Neurons emit a discrete number of spikes, but their responses often need to represent continuous quantities. Most studies of neural sampling implicitly rely on approximating Poissonian spike counts with Gaussian variables (e.g., refs. 29,31,51). However, this approximation does not work well when only a few spikes are emitted. Here, we showed that discrete Poisson spike generation can be used to generate samples from a posterior distribution of a continuous stimulus feature using a temporally averaged, smooth population firing rate profile. Thus, we have shown how a sample from a continuous variable can be generated even with only a few spikes from the neuronal population. Moreover, conventional SBCs are used to generate samples directly in a neural space whose dimension is given by the number of neurons in the population<sup>16,27,28,30-34,50</sup>, where a neuronal response,  $\mathbf{r}_t$  is interpreted directly as a sample from the (marginal) posterior of neuronal responses,  $p(\mathbf{r})$ . Hence the posterior mean is the temporally averaged population response, and the covariance of population responses is the posterior covariance. In contrast, our proposed network generates samples in a low dimensional stimulus subspace embedded in high dimensional neural activity space. The linear projection of network activity,  $\mathbf{r}_t$ , onto the stimulus subspace represents a sample from the stimulus posterior, similar to a previous study<sup>29</sup>. A computational benefit of sampling in a low dimensional stimulus subspace is convergence speed, as the volume of the stimulus subspace is significantly smaller than that of the neural activity space. Indeed, in our examples sequences of samples generated by a single recurrent network (Fig. 4) and coupled networks (Fig. 6) can both converge to an equilibrium distribution in less than 20 ms, which is fast enough to complete inference on a behaviorally relevant time scale (Fig. S6). Furthermore, the multiplication of probability distributions of latent stimulus, which is central to Bayesian inference (e.g., cue combination, decision making, see review in ref. 15), can be implemented by summing the inputs to a neuronal population (Eq. (5)). This follows from the fact that the instantaneous population input (or firing rate) linearly encodes the logarithm of a probability distribution (Eqs. (1) and (5)). In contrast, producing samples in neural activity space using conventional SBCs requires nonlinear operations in neural circuits in order to multiply probability distributions (or histograms) of the samples<sup>15</sup>.

A recent study demonstrated that an E-I recurrent network of ratebased neurons can be numerically optimized for sampling-based Bayesian inference<sup>32</sup>. In contrast, we used a theoretical approach to derive a network model of simplified spiking neurons, which implements sampling-based inference. This allowed us to explicitly describe the putative neural mechanisms needed for such sampling. Although the two studies use different generative models and neural representations, the network models in both studies share some common characteristics: ring structure, Poisson-like response variability, and tuning-dependent noise correlation (Fig. 1d). This implies that the seemingly different generative models and neural representations in the two studies reflect more general principles, as suggested in<sup>51</sup>. It will be interesting to extend our theoretical approach to dynamical spiking neurons to determine how the timescales of neuronal dynamics and neuronal oscillations impact inference in rich, dynamic sensory scenes (see below).

Differential noise correlations generated by recurrent network interactions are a signature of network sampling in our framework (Figs. 5c and 8c). This is in contrast to earlier studies where differential correlations were inherited from feedforward inputs<sup>17,52</sup>. While internally generated differential correlations could also emerge from a recurrent circuit which is not implementing inference<sup>22,24,52-55</sup> or implementing inference via other algorithms<sup>56</sup>, in our framework, the relation between the magnitude of internally generated differential correlations, the posterior uncertainty, and the strength of the recurrent synaptic weight (Eq. (9)) provides a clear test which can be used to verify our proposed circuit mechanism of sampling-based inference. One possible experimental approach would modulate the functional recurrent strength by using a perceptual learning task. Specifically, after using a reference stimulus set with a prescribed correlation between latent stimuli to fully train an animal, we expect that recurrent synaptic weights will strengthen or weaken to improve inference (Fig. 8e, dashed line). This will result in a fixed value of differential noise correlations in the population response due to the recurrent circuitry. Re-training with a stimulus set that has more (less) correlated latent stimuli compared to the reference set will cause the recurrent weights to increase (decrease) (Fig. 8e, red line). When the reference stimulus set is again used to drive task behavior, then performance (as a proxy of mutual information) will decrease, regardless of whether differential correlations have increased or decreased compared to those resulting from the reference stimulus set (Fig. 8e, arrows). In brief, the non-monotonic relationship between differential noise correlations and the mutual information between stimulus and responses which support Bayesian inference offers a clear (and falsifiable) experimental prediction.

Implementing sampling-based inference in our proposed network requires that feedforward and recurrent inputs have the same tuning

profile over the stimulus (Eq. (5)). This assumption is supported by experiments in layers 4 and 2/3 in mouse V18. Moreover, the recurrent connections in our network model are translation-invariant in the stimulus subspace, an assumption widely used in studies of continuous attractor networks (CAN)22,54,57,58, and a recent network model implementing sampling<sup>32</sup>. Perfectly translation-invariant connections are not strictly required for a circuit to implement sampling, but this assumption allows us to simplify the mathematical analysis. Adding randomness in recurrent connectivity would increases the variance of sampling distributions. We could then adjust the overall recurrent weight (a scalar) so that the sampling distribution matches the posterior, with no need to fine-tune individual synaptic weights in the network model. In the past, CANs have been shown to achieve maximal likelihood estimation (point estimate) via template matching<sup>15,58,59</sup>. Here we have shown that a network with CAN-like structure and internally Poisson spiking variability is able to perform sampling-based Bayesian inference. In our network correlations in the stimulus prior are represented by the strength of recurrent synaptic activity, which implies that the (subjective) prior precision in the network increases with the feedforward input strength.

To maintain a fixed prior in the network recurrent weights need to decrease with increased feedforward input strength which encodes the likelihood precision,  $\Lambda_f$  (Eq. (6)). Therefore, the (subjective) prior stored in the network with fixed recurrent weights may differ from the objective stimulus prior in the world ( $\Lambda_s$  in Eqs. (3) and (7)) with feedforward inputs of different strengths. One possibility is that the proposed network model does not generate samples from each distinct posterior determined by a specific feedforward input,  $p(\mathbf{s}|\mathbf{u}^{t})$ , but rather generates samples from the average sampling distribution over all possible feedforward inputs and hence matches the average posterior distribution  $\mathbb{E}_{p(\mathbf{u}^f)}[p(\mathbf{s}|\mathbf{u}^f)] = \mathbb{E}_{p(\mathbf{u}^f)}[q(\mathbf{s}|\mathbf{u}^f)],$ where  $\mathbb{E}_{p(\mathbf{u}^f)}[\cdot]$  denotes the average over the distribution  $p(\mathbf{u}^f)$ . Since the proposed recurrent circuit is general, this result may explain one source of inductive bias in cortical processing<sup>60</sup>. On the other hand, sampling correctly from each specific posterior could be achieved using different biophysical mechanisms that can modulate synaptic strengths and that we have not included in our model. For instance, short-term synaptic depression<sup>61</sup> or spike frequency adaptation<sup>62</sup> are gain control mechanisms that would allow the recurrent input strength (representing the prior correlation) to remain relatively fixed despite an increase in the feedforward input strength. Another possibility is that the recurrent circuit represents a more complex generative model which better captures the statistical structure of natural stimuli<sup>30,32,63</sup>. Here we assumed that the generative models represented by the network match the model that generate the sensory stimuli. This is unlikely to be the case in practice. Such mismatch between the true and internal model of the world can lead to biases and increased noise which are likely to manifest in specific ways in neural circuits that perform inference via sampling<sup>64</sup>. Furthermore, we only considered sampling driven by spiking variability with a Fano factor of 1, while cortical responses often have Fano factors that differ from 165,66. In the latter case, our theory can still work by changing the feedforward connection weight to compensate for the change in Fano factor, as suggested in a recent study<sup>67</sup>.

To keep our exposition transparent, we only presented models with minimal complexity. Our proposed network mechanism of sampling-based inference can be generalized to more complex generative models, since the assumption of Gaussianity (Eqs. (21) and (22)) and the analytical expression in Eq. (24) are not essential, and several relaxed frameworks may be explored. First, similar networks can generate samples from other multi-dimensional distributions where the conditional distribution of each latent variable belongs to the linear exponential family<sup>38,39</sup>. This could be done by changing the tuning functions of neurons to another appropriate profile, as the logarithm of tuning determines the type of sampling distribution (Eq.

(1)). When sampling from non-Gaussian distributions, the stimulus samples can be linearly read out with the weight determined by the tuning profile (i.e.,  $\mathbf{h}(s)$  in Eq. (1)<sup>39</sup>,). Second, the tuning of recurrent inputs does not need to be the same as that of feedforward inputs. Instead, the logarithm of recurrent input tuning can have a form of the conjugate prior with the likelihood conveyed by feedforward inputs. Third, the network model could also be used to infer the latent variables with a non-uniform marginal prior, if, for example, the preferred stimuli of neurons in the population are not distributed uniformly in the stimulus subspace<sup>68</sup>. And the proposed network model has the potential to produce samples from the posterior distribution of latent dynamic stimuli which can be described by a hidden Markov model. Lastly, we considered only non-structured inhibition for simplicity. Structured inhibitory connections could modulate the position of excitatory responses in the stimulus subspace, i.e., the mean of the conditional distribution. Such interplay between E and I neurons with structured inhibition has the potential to implement Hamiltonian sampling, where the I neurons represent the sample of auxiliary variables<sup>38,50</sup>.

In conclusion, we have shown that a recurrent circuit of neurons with Poisson spiking statistics can implement sampling from a family of multivariate posterior distributions, with internal spiking variability driving the generation of stimulus samples, and the recurrent connections representing the stimulus prior. The proposed neural code may help us understand the structure of neuronal activity, and provide a building block for more complicated population computations.

#### Methods

#### A linear network of excitatory neurons

We study how a generic recurrent network model consisting solely of  $N_E$  excitatory (E) neurons with Poisson spiking statistics (no inhibitory neurons) can implement sampling-based Bayesian inference to approximate the stimulus posterior. We describe neuronal activity using a time-discretized Hawkes process (a type of multivariate, inhomogeneous Poisson process<sup>69</sup>). The instantaneous firing rates of the neurons in the network at time t,  $\lambda_t$ , obey the following recurrent equations:

$$\lambda_t \Delta t = \mathbf{u}^{\mathsf{f}} + \mathbf{u}_t^{\mathsf{r}} = \mathbf{u}^{\mathsf{f}} + (w_E \mathbf{r}_{t-\Delta t} + \sigma_r \boldsymbol{\xi}_t), \tag{10}$$

$$\mathbf{r}_{t} \sim \prod_{i=1}^{N_{E}} \operatorname{Poisson}\left(\boldsymbol{\lambda}_{tj}\Delta t\right),$$
 (11)

where  $\mathbf{u}^t$  is the feedforward Poisson spiking input (described below; Eq. (18)),  $\mathbf{u}_t^r$  is the continuous valued recurrent input at time t, and  $\boldsymbol{\xi}_t$  is a  $N_E$  dimensional independent Gaussian white noise. Hence, over each time interval  $[t - \Delta t, t]$  the activity of the neurons in the network is modeled by a vector of independently generated Poisson spike counts,  $\mathbf{r}_t$ , with means determined by the rates  $\lambda_t$ . The parameters  $w_E$  and  $\sigma_t$  determine the excitatory recurrent weight and recurrent variability, respectively. The instantaneous firing rate  $\lambda_t$  can be negative due to the recurrent input and noise (Eq. (36)). We interpret a negative firing rate,  $\lambda_t$ , as a zero probability of generating a spike.

#### Poisson spike generation samples stimulus

Independent Poisson spike generation in the network whose activity is described by Eq. (11) can drive sampling across time or across trials from a conditional stimulus distribution determined by the instantaneous firing rate  $\lambda_t$ . Below, we compute the distribution of stimulus samples given  $\lambda_t$ . We assume that the instantaneous firing rate,  $\lambda_t$ , has a smooth bell-shaped profile and can be parameterized as,

$$\lambda_{ti} = R \exp[-(\bar{s}_t - \theta_i)^2 / 2a^2] = R \exp[\mathbf{h}_i(\bar{s}_t)], \tag{12}$$

where  $\bar{s}_t$  characterizes the position of the population firing rate on the stimulus subspace (Fig. 1b, x-axis), while R and a denote the height and width of the population firing rate, respectively. Further,  $\theta_j$  is the preferred stimulus value of neuron j, and the preferred stimuli of all neurons,  $\{\theta_j\}_{j=1}^{N_E}$ , are uniformly distributed over the range of stimulus s (Fig. 1b).

To simplify the analysis, we first assume that the instantaneous firing rate is fixed over time. When generating Poisson spikes  $\mathbf{r}_t$  from  $\lambda_t$ , the probability of observing a stimulus sample  $\tilde{s}_t$  (embedded in  $\mathbf{r}_t$ ) can be derived as (see details in Supplementary Information),

$$p(\mathbf{r}_{t}|\boldsymbol{\lambda}_{t}) = \prod_{j=1}^{N_{E}} \operatorname{Poisson}\left(\mathbf{r}_{tj}|\boldsymbol{\lambda}_{tj}\Delta t\right),$$

$$\propto \exp[\mathbf{h}(\bar{s}_{t})^{\top}\mathbf{r}] \cdot \left[n_{A}^{n_{r}} \exp(-n_{A})\right],$$

$$\propto \mathcal{N}\left(\tilde{s}_{t}|\bar{s}_{t}, a^{2}n_{r}^{-1}\right) \operatorname{Poisson}(n_{r}|n_{A}),$$
(13)

where  $n_{\bf r}=\Sigma_j{\bf r}_{tj}$  is the number of emitted spikes across the whole neural population, and  $n_{\lambda}=\Sigma_j\langle\lambda_j\rangle\Delta t$  is the sum of population firing rate. Here  $\mathcal{N}(s|\mu,\sigma^2)$  denotes a Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ , and  ${\bf h}(\bar{s}_t)$  is a vector with the jth element as  ${\bf h}_j(\bar{s}_t)$  shown in Eq. (12). The logarithm of the firing rate profile,  ${\bf h}(\bar{s}_t)$ , determines how the stimulus sample  $\bar{s}_t$  and its mean,  $\bar{s}_t$ , can be read out respectively from  ${\bf r}_t$  and  $\lambda_t$ .

$$\tilde{s}_t = \sum_j \mathbf{r}_{tj} \theta_j / \sum_j \mathbf{r}_{tj}, \quad \bar{s}_t = \sum_j \lambda_{tj} \theta_j / \sum_j \lambda_{tj},$$
(14)

where  $\tilde{s}_t$  and  $\bar{s}_t$  characterizes the position of  $\mathbf{r}_t$  and  $\boldsymbol{\lambda}_t$  on the stimulus subspace.

The sampling variability of  $\tilde{s}_t$  in a single time step depends on the number of emitted spikes,  $n_{\rm r}$ . When the fixed rates,  $\lambda_t$ , repeatedly generate spikes over time, the sampling distribution of  $\tilde{s}_t$  can be calculated by marginalizing the likelihood (Eq. (13), last line) over different values of  $n_{\rm r}$  since  $n_{\rm r}$  varies across time (detailed calculation by using Laplacian approximation can be seen in Supplementary Information),

$$p(\tilde{s}_{t}|\boldsymbol{\lambda}_{t}) = \sum_{n_{r}} \mathcal{N}(\tilde{s}_{t}|\bar{s}_{t},a^{2}n_{r}^{-1}) \operatorname{Poisson}(n_{r}|n_{\lambda}),$$

$$\approx \mathcal{N}(\tilde{s}_{t}|\bar{s}_{t},a^{2}n_{\lambda}^{-1}).$$
(15)

Each stimulus sample,  $\tilde{s}_t$ , is thus drawn from a conditional distribution determined by the instantaneous firing rate,  $p(\tilde{s}|\boldsymbol{\lambda}_t)$ , and can be written as

$$\tilde{s}_t \sim p(\tilde{s}|\boldsymbol{\lambda}_t) = \mathcal{N}\left(\tilde{s}|\bar{s}_t, a^2 n_{\boldsymbol{\lambda}}^{-1}\right) \propto \exp[\boldsymbol{h}(\tilde{s})^\top \boldsymbol{\lambda}_t]. \tag{16}$$

The last proportionality in the above equation is satisfied by a Gaussian profile in the firing rate (more general derivation can be found in Supplementary Information). Introducing  $\Lambda = a^{-2}n_{\lambda}$  gives Eq. (1) shown in the main text.

Eq. (16) suggests that the type of sampling distribution (or the conditional distribution) that is obtained from spike generation variability is determined by the profile of the instantaneous firing rate, i.e.,  $\mathbf{h}(\bar{s}_t)$  (Eq. (12)). Although the sampling distribution belongs to the linear exponential family of distributions which is similar to the probabilistic population code (PPC)<sup>39</sup>, there are different ways in representing these distributions. In PPCs, the likelihood over  $\bar{s}_t$  is parametrically represented by a single realization of independent neuronal response  $\mathbf{r}$  (Eq. (13)), while in our work the distribution is approximated by a sequence of samples,  $\tilde{s}_t$ , effectively generated by conditionally independent Poisson spike discharges.

The above analysis can be extended to the case where the instantaneous firing rate,  $\lambda_t$ , in a time step deviates from a smooth Gaussian profile (Eq. (12)), which is the case in the actual network simulations. In general,  $\lambda_t$  can be expressed as,

$$\boldsymbol{\lambda}_{tj} = R_t \exp[\mathbf{h}_j(\bar{s}_t)] + \boldsymbol{\delta}_{\perp} \boldsymbol{\lambda}_{tj}, \tag{17}$$

where  $\delta_{\perp} \lambda_t$  denotes the deviation from a smooth Gaussian profile. Note that the sampling distribution only depends on the position,  $\bar{s}_t$ , and the sum of instantaneous firing rate,  $n_{\lambda}$  (Eq. (16)), which corresponds to two perpendicular directions in the  $N_E$  dimensional space of  $\lambda_t$ . For any instantaneous firing rate vector,  $\lambda_t$ , we can always find  $\bar{s}_t$  and  $R_t$  that make the deviation  $\delta_{\perp} \lambda_t$  perpendicular to the two directions, i.e.,  $\sum_j \delta_{\perp} \lambda_{tj} \theta_j = 0$ , and  $\sum_j \delta_{\perp} \lambda_{tj} = 0$ . This observation implies that deviations from Gaussian firing rate profiles do not affect our theory.

## Feedforward spiking input conveys the likelihood of stimulus

We model the feedforward inputs to the E neurons in the network,  $\mathbf{u}^{t}$ , as independent Poisson spikes, with Gaussian tuning over stimulus s,

$$p(\mathbf{u}^{f}|s) = \prod_{j=1}^{N_{E}} \text{Poisson}\left[\mathbf{u}_{j}^{f}|\langle \mathbf{u}_{j}^{f}(s)\rangle\right],$$

$$\langle \mathbf{u}_{j}^{f}(s)\rangle = U^{f} \exp[\mathbf{h}_{j}(s)] = U^{f} \exp[-(\theta_{j} - s)^{2}/2a^{2}].$$
(18)

Here  $\mathbf{u}_{j}^{f}$  denotes the feedforward input received by the jth E neuron, and  $\langle \mathbf{u}_{j}^{f}(s) \rangle$  is the tuning of the feedforward input. This mathematical description of feedforward input is the same as the one used in the definition of typical PPCs<sup>15,39,40</sup>. Since the preferred stimulus values,  $\{\theta_{j}\}_{j=1}^{N_{E}}$ , of all feedforward inputs are uniformly distributed in stimulus space then the likelihood of s given a single observation of the input,  $\mathbf{u}^{f}$ , satisfies<sup>39,40</sup>,

$$p(\mathbf{u}^{\mathsf{f}}|s) \propto \exp\left[\mathbf{h}(s)^{\mathsf{T}}\mathbf{u}^{\mathsf{f}}\right],$$

$$\propto \mathcal{N}\left(s|\mu_{\mathsf{f}},\Lambda_{\mathsf{f}}^{-1}\right).$$
(19)

The logarithm of tuning,  $\mathbf{h}(s)$ , determines the type of likelihood<sup>15</sup>. Specifically, the Gaussian tuning leads to a Gaussian likelihood (Eq. (19)), whose mean,  $\mu_t$ , and precision,  $\Lambda_t$ , are both linear functions of the inputs,

$$\mu_{\rm f} = n_{\rm f}^{-1} \sum_{i} \mathbf{u}_{j}^{\rm f} \theta_{j}, \quad \Lambda_{\rm f} = a^{-2} n_{\rm f} = a^{-2} \sum_{i} \mathbf{u}_{j}^{\rm f}.$$
 (20)

The mean,  $\mu_t$ , represents the position of  $\mathbf{u}^t$  in stimulus subspace, and the precision,  $\Lambda_t$ , is proportional to the sum of total feedforward spike counts,  $n_t$ .

# A recurrent network samples hierarchical latent variables

A hierarchical generative model. We consider a hierarchical generative model for which inference can be implemented in a recurrent circuit of Poisson neurons. We extend the simple generative model of feedforward input (Eq. (19)) by considering the stimulus s to depend on a one-dimensional stimulus parameter variable, z. For simplicity, we assume that z follows a uniform distribution (Fig. 3b, marginal plots)

$$p(z) = \mathcal{U}(-180^{\circ}, 180^{\circ}),$$
 (21)

where  $\mathcal{U}(a,b)$  denotes a uniform distribution over [a,b]. The assumption of a uniform prior, p(z), simplifies our model significantly, as it implies the spatial homogeneity of the network model as given by Eqs. ((18), (19)). However, this assumption is not essential for our main results. Due to the differences between the stimulus and its underlying

parameter of the sensory scene, the stimulus, *s*, is not identical to the parameter *z*, but we assume that the two are correlated, so that

$$p(s|z,\Lambda_s) = \mathcal{N}\left(s|z,\Lambda_s^{-1}\right). \tag{22}$$

In sum, the whole generative model is determined by,

$$p(\mathbf{u}^{\mathsf{f}}, s, z) = p(\mathbf{u}^{\mathsf{f}}|s)p(s|z)p(z),$$

$$\propto \mathcal{N}\left(s|\mu_{\mathsf{f}}, \Lambda_{\mathsf{f}}^{-1}\right) \mathcal{N}\left(s|z, \Lambda_{\mathsf{s}}^{-1}\right),$$
(23)

where  $p(\mathbf{u}^{t}|s)$  is the same as in Eq. (19).

Approximate Bayesian inference via Gibbs sampling. The joint posterior of s and z can be analytically derived given the generative model (Eq. (23)),

$$p(s,z|\mathbf{u}^{\mathsf{f}}) = \mathcal{N}\left[(s,z)^{\top}|\boldsymbol{\mu}_{p},\mathbf{K}_{p}^{-1}\right],$$

$$\boldsymbol{\mu}_{p} = (\boldsymbol{\mu}_{\mathsf{f}},\boldsymbol{\mu}_{\mathsf{f}})^{\top}, \quad \mathbf{K}_{p} = \begin{pmatrix} \Lambda_{\mathsf{f}} + \Lambda_{s} & -\Lambda_{s} \\ -\Lambda_{s} & \Lambda_{s} \end{pmatrix}.$$
(24)

We will use this expression to verify that the samples produced by our algorithm converge to the output of the algorithm.

We use the stochastic response of our recurrent network (Eqs. (10), (11)), as a basis for Gibbs sampling<sup>31,38,42</sup> (a type of Monte Carlo method) to approximate the joint posterior p(s,z). To describe the iterative Gibbs algorithm, we assume that a stimulus parameter sample,  $\tilde{z}_t$ , is provided at time t, which is then combined with the feedforward input to update the conditional distribution of stimulus s (step 1 in Fig. 3c),

$$p(\tilde{s}|\tilde{z}_{t}, \mathbf{u}^{f}) \propto p(\mathbf{u}^{f}|\tilde{s})p(\tilde{s}|\tilde{z}_{t}) \propto \mathcal{N}\left(s|\bar{s}_{t}, \Lambda^{-1}\right),$$

$$\bar{s}_{t} = \frac{\Lambda_{f}\mu_{f} + \Lambda_{s}\tilde{z}_{t}}{\Lambda_{f} + \Lambda_{s}}, \quad \Lambda = \Lambda_{f} + \Lambda_{s}.$$
(25)

The next step in the algorithm is to draw a sample,  $\tilde{s}_t$ , from the conditional distribution  $p(\tilde{s}|\tilde{z}_t, \mathbf{u}^f)$  (step 2 in Fig. 3c),

$$\tilde{s}_t \sim p(\tilde{s}|\tilde{z}_t, \mathbf{u}^{\mathsf{f}}) = \mathcal{N}\left(\tilde{s}|\bar{s}_t, \Lambda^{-1}\right).$$

Next, the conditional distribution of stimulus parameter, z, is updated given this new sample,  $\bar{s}_t$ , and a new sample,  $\tilde{z}_{t+\Delta t}$ , is drawn (step 3 in Fig. 3c),

$$\tilde{z}_{t+\Lambda t} \sim p(\tilde{z}|\tilde{s}_t) = \mathcal{N}(\tilde{z}|\tilde{s}_t, \Lambda_s^{-1}).$$
 (26)

These three steps in the Gibbs sampling algorithm (Eqs. (25), (26)) are performed iteratively until sufficiently many samples,  $\tilde{s}_t$  and  $\tilde{z}_t$ , are generated to approximate the true posterior distribution with sufficient accuracy (Fig. 3d; compare the red dots with the blue contour map).

Implementing the Gibbs sampling in a recurrent circuit model. Gibbs sampling of the stimulus (Eq. (4b)) can be implemented via independent Poisson spike generation, as long as the conditional distribution encoded in  $\lambda_t$  (Eq. (16)) is the same as the conditional distribution in the Gibbs sampling algorithm (Eq. (4a)), i.e.,  $\ln p(\tilde{s}|\lambda_t) = \mathbf{h}(\tilde{s})^{\top} \lambda_t = \ln p(\tilde{s}|\tilde{z}_t, \mathbf{u}^t)$ . This condition can be realized in the recurrent circuit by relating the expressions describing the neural dynamics (Eq. (10)) and those describing the Gibbs sampling

distribution (Eq. (4a)) to yield,

$$\ln p(\tilde{\mathbf{s}}|\tilde{\mathbf{z}}_{t},\mathbf{u}^{f}) = \mathbf{h}(\tilde{\mathbf{s}})^{\top} \boldsymbol{\lambda}_{t},$$

$$= \mathbf{h}(\tilde{\mathbf{s}})^{\top} \mathbf{u}^{f} + \mathbf{h}(\tilde{\mathbf{s}})^{\top} \mathbf{u}^{r}_{t},$$

$$= \ln p(\mathbf{u}^{f}|\tilde{\mathbf{s}}) + \ln p(\tilde{\mathbf{s}}|\tilde{\mathbf{z}}_{t}).$$
(27)

The generative model for the feedforward input  $\mathbf{u}^{f}$  (Eq. (19)) suggests that  $\ln p(\mathbf{u}^{f}|\tilde{\mathbf{s}}) = \mathbf{h}(\tilde{\mathbf{s}})^{\top}\mathbf{u}^{f}$ . Hence to satisfy Eq. (27) we require

$$\ln p(\tilde{s}|\tilde{z}_t) = \mathbf{h}(\tilde{s})^{\top} \mathbf{u}_{t,t}^{r}$$
(28)

which implies that the recurrent input,  $\mathbf{u}_t^r$ , should approximately have a Gaussian profile,

$$\mathbf{u}_{tj}^{\mathrm{r}}(\tilde{z}_{t}) = U^{\mathrm{r}} \exp[-(\theta_{j} - \tilde{z}_{t})^{2}/2a^{2}] + \delta_{\perp} \mathbf{u}_{tj}^{\mathrm{r}},$$

$$\tilde{z}_{t} = \sum_{i} \mathbf{u}_{tj}^{\mathrm{r}} \theta_{j} / \sum_{i} \mathbf{u}_{tj}^{\mathrm{r}}, \quad \Lambda_{s} = a^{-2} \sum_{i} \mathbf{u}_{tj}^{\mathrm{r}},$$
(29)

whose position on the stimulus subspace is  $\tilde{z}_t$ , and the sum of input (height) is determined by  $\Lambda_s$ , the precision of conditional distribution  $p(s|\tilde{z}_t)$ . In a similar fashion to Eq. (17),  $\delta_{\perp} \mathbf{u}_t^r$  denotes the deviation from a smooth Gaussian and is perpendicular to the direction of  $\tilde{z}_t$  and  $\Lambda_s$ .

The optimal recurrent weight can be derived by combining Eq. (29) and Eq. (17). We notice the recurrent input,  $\mathbf{u}'$ , and neuronal responses,  $\mathbf{r}_t$ , have the same tuning width, a, in a network with only E neurons. This can only be achieved if E neurons are only self-connected (Eq. (10)), as lateral connection broaden their tuning. The optimal recurrent weight generating recurrent input with appropriate strength is then,

$$w_E^* = \frac{\langle \mathbf{u}_j^r \rangle}{\langle \mathbf{r}_j \rangle} = \frac{\sum_j \langle \mathbf{u}_j^r \rangle}{\sum_j \langle \mathbf{r}_j \rangle} = \frac{\sum_j \langle \mathbf{u}_j^r \rangle}{\sum_j \left( \langle \mathbf{u}_j^r \rangle + \langle \mathbf{u}_j^r \rangle \right)} = \frac{\Lambda_s}{\Lambda_f + \Lambda_s},$$
 (30)

which yields Eq. (6) in the main text. Note that the self-connection is a result of the simplifying assumption that the network consists solely of E neurons (Eq. (10)), which can be relaxed in a full network consisting both E and I neurons as we show below.

The sampling of the stimulus parameter (Eq. (4c)) can be implemented through variability in the recurrent input. To do this, we include diffusive term in the recurrent interactions,  $\mathbf{u}_t^r$ , and we equate the variance of the fluctuations with the mean to mimic a Poisson distribution:

$$\mathbf{u}_{t}^{\mathsf{r}} = \bar{\mathbf{u}}_{t}^{\mathsf{r}} + \sqrt{[\bar{\mathbf{u}}_{t}^{\mathsf{r}}]_{+}} \boldsymbol{\xi}_{t}, \quad \bar{\mathbf{u}}_{t}^{\mathsf{r}} = \boldsymbol{w}_{E}^{*} \mathbf{r}_{t-\Delta t}, \tag{31}$$

where  $[\cdot]_+$  denotes negative rectification. Here  $\xi_t$  is a  $N_E$  dimensional Gaussian white noise with  $\langle \xi_t(i)\xi_{t'}(j)\rangle = \delta_{ij}\delta(t-t')$ ,  $\delta_{ij}$  and  $\delta(t-t')$  are Kronecker and Dirac delta functions respectively,  $\bar{\mathbf{u}}_t^r$  represents the conditional distribution  $p(\tilde{z}|\tilde{s}_{t-\Delta t})$ , and  $\mathbf{u}_t^r$  represent a stimulus parameter sample  $\tilde{z}_t$  (Eq. (29)). The multiplicative variability on recurrent interaction may come from synaptic noise<sup>37,70</sup>.

## Coupled circuits sample a multi-dimensional posterior

We consider a generative model which has multiple latent stimuli,  $\mathbf{s} = (s_1, s_2, \dots, s_m)$ , which are organized in parallel (Fig. 6a). Without loss of generality, we consider the simplest case where m = 2, and the same mechanism can be straightforwardly extended to any m > 2. We assume the joint prior of  $\mathbf{s}$  is a multivariate normal distribution,

$$p(\mathbf{s}) = \mathcal{N}(\mathbf{s}|\boldsymbol{\mu}_{s},\boldsymbol{\Lambda}_{s}^{-1}) \propto \exp[-\Lambda_{s}(s_{1} - s_{2})^{2}/2],$$
with  $\boldsymbol{\Lambda}_{s} = \Lambda_{s} \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}$ , (32)

and each stimulus  $s_m$  is uniformly distributed in (-180°, 180°) with periodic boundary imposed. The definition of Gaussian distribution in a circular space works well as long as the variance of the distribution is much smaller than the range of stimulus space. Here  $\Lambda_s$  is the precision matrix, while the scalar variable  $\Lambda_s$  ( $\Lambda_s \ge 0$ ) characterizes the correlation between  $s_1$  and  $s_2$ . Note that the covariance matrix  $\mathbf{\Lambda}_s^{-1}$  is not defined, and the prior (Eq. (32)) is improper. The mean,  $\mu_s$ , is a free parameter, because it doesn't appear in the detailed expression of the prior (Eq. (32)), which is a consequence from the zero determinant of the precision matrix, i.e.,  $|\mathbf{\Lambda}_s| = 0$ . A further consequence is that the prior is not centered at  $\mu_s$ , but instead has a band structure along the diagonal, and the marginal prior of each stimulus feature  $p(s_m)$  (m=1,2) is uniform (Fig. 6b). The uniform marginal prior simplifies our theoretical derivation as it implies the spatial homogeneity of the network model but doesn't impact the proposed neural coding

Each stimulus  $s_m$  (m=1,2) individually generates feedforward spiking input  $\mathbf{u}_m^f$ , whose likelihood  $p(\mathbf{u}_m^f|s_m)$  is exactly the same as Eq. (2). Combined together, the generative model is

$$p(\mathbf{u}^{\mathsf{f}}|\mathbf{s})p(\mathbf{s}) = \left[\prod_{m=1}^{2} p(\mathbf{u}_{m}^{\mathsf{f}}|s_{m})\right] p(s_{1},s_{2}),$$

$$\propto \left[\prod_{m=1}^{2} \mathcal{N}(s_{m}|\mu_{\mathsf{f}m},\Lambda_{\mathsf{f}m}^{-1})\right] \mathcal{N}(\mathbf{s}|\boldsymbol{\mu}_{\mathsf{s}},\boldsymbol{\Lambda}_{\mathsf{s}}^{-1}),$$

$$\propto \mathcal{N}(\mathbf{s}|\boldsymbol{\mu}_{\mathsf{f}},\boldsymbol{\Lambda}_{\mathsf{f}}^{-1})\mathcal{N}(\mathbf{s}|\boldsymbol{\mu}_{\mathsf{s}},\boldsymbol{\Lambda}_{\mathsf{s}}^{-1}),$$
(33)

where  $\mu_f = (\mu_{fl}, \mu_{f2})^T$ , and the likelihood precision matrix  $\Lambda_f = \text{diag}(\Lambda_{fl}, \Lambda_{f2})$  is a diagonal matrix.

Gibbs sampling of the multi-dimensional posterior in a coupled neural circuit. Given the generative model (Eq. (33)), the joint posterior of  $s_1$  and  $s_2$  is a bivariate normal distribution, i.e.,  $p(\mathbf{s}|\mathbf{u}^f) = \mathcal{N}\left(\mathbf{s}|\mathbf{u}_p, \mathbf{K}_p^{-1}\right)$ , whose precision matrix  $\mathbf{K}_p$  and the mean  $\mathbf{\mu}_p$  are,

$$\mathbf{K}_{n} = \boldsymbol{\Lambda}_{f} + \boldsymbol{\Lambda}_{S}, \quad \boldsymbol{\mu}_{n} = \mathbf{K}_{n}^{-1} \boldsymbol{\Lambda}_{f} \boldsymbol{\mu}_{f}. \tag{34}$$

The precision matrix of the posterior is the sum of the precision of the likelihood and the prior, implying increased reliability of the distribution after combining with the prior. Meanwhile, the posterior mean is the weighted average of the means of the two likelihoods, with the weight proportional to the precision of each likelihood. We use this expression for the posterior to evaluate the performance of the proposed sampling-based algorithm.

Using Gibbs sampling to approximate the posterior (Eq. (34)) involves the following steps:

Compute: 
$$p(\tilde{s}_1|\mathbf{u}_1^f, \tilde{s}_{2,t-\Delta t}) \propto p(\mathbf{u}_1^f|\tilde{s}_1)p(\tilde{s}_{2,t-\Delta t}|\tilde{s}_1)$$
, (35a)

Sample: 
$$\tilde{\mathbf{s}}_{1t} \sim p(\tilde{\mathbf{s}}_1 | \mathbf{u}_{1}^{\mathsf{f}}, \tilde{\mathbf{s}}_{2t-\Lambda t})$$
. (35b)

We note that we only describe the sampling from the posterior distribution of  $s_1$ ; as samples from the posterior of  $s_2$  can be obtained similarly after exchanging indices. This sampling can be implemented in a neural circuit model consisting of several coupled networks, in which each network generates samples from the posterior distribution of the corresponding stimulus. Therefore, the number of networks in the coupled circuit equals the dimension of the latent stimuli. The dynamics of the coupled neural circuit is defined by:

$$\lambda_{1t} = \mathbf{u}_1^{\mathsf{f}} + \mathbf{u}_{12\,t}^{\mathsf{r}} = \mathbf{u}_1^{\mathsf{f}} + w_{12}\mathbf{r}_{2\,t-\Lambda t},\tag{36}$$

$$\mathbf{r}_{1t} \sim \prod_{j=1}^{N_E} \text{Poisson}(\boldsymbol{\lambda}_{1t,j}),$$
 (37)

We again note the dynamics of network 2 can be similarly obtained by changing indices. To implement Gibbs sampling (Eqs. (35a), (35b)) in the coupled circuit (Eqs. (36), (37)), spike generation in network 1 (Eq. (37)) can be used to produce stimulus samples,  $\tilde{s}_{1t}$ , when the conditional distribution determined by  $\lambda_{1t}$  matches the conditional distribution required in the definition of Gibbs sampling (Eq. (35a)), i.e.,  $\ln p(\tilde{s}_1|\mathbf{u}_1^t,\tilde{s}_{2,t-\Delta t}) = \ln p(\tilde{s}_{1t}|\lambda_{1t}) = \mathbf{h}(\tilde{s}_1)^\top \lambda_{1t}$ . Taking the logarithm of Eq. (35a) yields,

$$\ln p(\tilde{s}_1|\mathbf{u}_1^f, \tilde{s}_{2,t-\Delta t}) = \ln p(\mathbf{u}_1^f|\tilde{s}_1) + \ln p(\tilde{s}_{2,t-\Delta t}|\tilde{s}_1). \tag{38}$$

Comparing this expression with Eq. (36), we see that the feedforward input,  $\mathbf{u}_1^f$ , matches the conditional distribution  $p(\mathbf{u}_1^f|\tilde{s}_1)$  (Eq. (33)). We therefore require the recurrent input from network 2 to network 1 to encode the conditional distribution  $p(\tilde{s}_{2,t-\Delta t}|\tilde{s}_1)$ , i.e.,  $\ln p(\tilde{s}_{2,t-\Delta t}|\tilde{s}_1) = \mathbf{h}(\tilde{s}_1)^{\mathsf{T}}\mathbf{u}_{12,t}^r$ . This implies that  $\mathbf{u}_{12,t}^r$  should approximately have a Gaussian profile,

$$\mathbf{u}_{12,tj}^{r} = U_{12}^{r} \exp[-(\theta_{j} - \tilde{s}_{t-\Delta t})^{2}/2a^{2}] + \delta_{\perp} \mathbf{u}_{12,tj}^{r}, \tilde{s}_{2,t-\Delta t} = \sum_{i} \mathbf{u}_{12,tj}^{r} \theta_{j} / \sum_{i} \mathbf{u}_{12,tj}^{r}, \quad \Lambda_{s} = a^{-2} \sum_{i} \mathbf{u}_{12,tj}^{r},$$
(39)

where  $\delta_{\perp} \mathbf{u}_{12,tj}^{r}$  quantifies the deviation from a perfect Gaussian profile, and does not affect the decoded value  $\tilde{s}_{2,t-\Delta t}$  and  $\Lambda_{s}$ .

The recurrent input,  $\mathbf{u}_{12}^r$ , (Eq. (39)) has the same width a as the neuronal response,  $\mathbf{r}_1$ . In circuit containing only E neurons, if the two networks have the same number of neurons, then across networks only neurons having the same preferred stimulus should be connected. The optimal recurrent weight between two networks is then

$$w_{mn} = \frac{\langle \mathbf{u}_{mn,j}^r \rangle}{\langle \mathbf{r}_{nj} \rangle} = \frac{\sum_{j} \langle \mathbf{u}_{mn,j}^r \rangle}{\sum_{j} \langle \mathbf{r}_{nj} \rangle} = \frac{\Lambda_s}{\Lambda_s + \Lambda_n^f}, (m \neq n)$$
 (40)

Since each network individually generate a stimulus sample, the sample of stimulus m can be locally read out from network m's responses even if the activities of two networks are correlated (Fig. 6a), which greatly simplifies readout. Furthermore, due to the population firing rate of each network has Gaussian profile, the stimulus sample  $\tilde{s}_{mt}$  can be linearly read out from  $\mathbf{r}_{mt}$  as

$$\tilde{s}_{mt} = \sum_{j} \theta_{j} \mathbf{r}_{mt,j} / \sum_{j} \mathbf{r}_{mt,j}. \tag{41}$$

We note that the circuit implementation of Gibbs sampling from a multi-dimensional posterior (Eq. (8a)) does not require the recurrent connections between E neurons within a network. This is due to the assumption that the marginal priors of each stimulus feature,  $p(s_m)$ , are uniform. For a non-uniform marginal prior  $p(s_m)$ , recurrent connections between E neurons within a network would be required for generating samples from a distribution that matches the true posterior.

## Inference from an information-theoretic point of view

The goal of the sampling algorithm is to approximate the posterior distribution of a latent variables,  $\Theta$ , given a feedforward input,  $\mathbf{u}^t$ . Specifically, the latent variables  $\Theta = \{s, z\}$  in the hierarchical generative model (Eq. (23)), or  $\Theta = \mathbf{s} = \{s_1, s_2\}$  in the generative model with breadth (Eq. (33)). When the sampling algorithm uses an internal model which does not match the structure of the generative model, the sampling

distribution  $q(\Theta|\mathbf{u}')$  will differ from the true posterior,  $p(\Theta|\mathbf{u}')$  (Eq. (24)). In this case the mutual information between the sampling distribution of the latent variables,  $\Theta$ , and  $\mathbf{u}'$  will be smaller than in the case when samples come from the true posterior,  $p(\Theta|\mathbf{u}')$ ,

$$I(\Theta, \mathbf{u}^{f}) = -\mathbb{E}_{p(\Theta)}[\log p(\Theta)] + \mathbb{E}_{p(\Theta, \mathbf{u}^{f})}[\log p(\Theta|\mathbf{u}^{f})]$$

$$\geq -\mathbb{E}_{p(\Theta)}[\log p(\Theta)] + \mathbb{E}_{p(\Theta, \mathbf{u}^{f})}[\log q(\Theta|\mathbf{u}^{f})] \equiv I_{q}(\Theta; \mathbf{u}^{f}),$$
(42)

It is straightforward to show that the difference between  $I(\Theta, \mathbf{u}')$  and  $I_q(\Theta, \mathbf{u}')$  is the Kullback–Leibler (KL) divergence between p and q, i.e.,  $D_{KL}[p||q] = I(\Theta, \mathbf{u}^f) - I_q(\Theta, \mathbf{u}^f) = \mathbb{E}_p(\ln p - \ln q) \ge 0$ . Equality in Eq. (42) holds only if the distribution q matches the true posterior p.

The mutual information  $I_q(\Theta; \mathbf{u}^t)$  can be computed analytically when the approximating distribution  $q(\Theta|\mathbf{u}^t) = \mathcal{N}(\Theta|\boldsymbol{\mu}_q, \mathbf{K}_q^{-1})$  is a bivariate normal (substituting Eqs. (23) and (24) into Eq. (42)),

$$I_q(\Theta; \mathbf{u}^{\mathrm{f}}) = \log L + \frac{1}{2} \left[ 1 + \log \frac{|\mathbf{K}_q|}{2\pi\Lambda_{\mathrm{s}}} - \mathrm{tr}(\mathbf{K}_q\mathbf{K}_p^{-1}) - (\boldsymbol{\mu}_p - \boldsymbol{\mu}_q)^{\top} \mathbf{K}_q(\boldsymbol{\mu}_p - \boldsymbol{\mu}_q) \right]. \tag{43}$$

Here  $L=360^\circ$  is the length of the stimulus feature subspace, while  $\mu_p$  and  $\mathbf{K}_p$  are the mean and the precision matrix of the posterior distribution (Eqs. (24) or (34)). When q matches the posterior distribution, p, we have,  $I(\Theta; \mathbf{u}^f) = \log L - \frac{1}{2}[1 + \log(2\pi\Lambda_s) - \log |\mathbf{K}_p|]$ .

# The neuronal response distribution conditioned on external stimulus

We compute the distribution of neuronal responses  ${\bf r}$  over time/trial in response to an external stimulus s, i.e.,  $p({\bf r}|s)$ , in order to find a neural signature of network sampling and compare it with experimental data. For a fixed external stimulus s, the neuronal response  ${\bf r}$  fluctuates due to both sensory transmission noise described by  $p({\bf u}^i|s)$  (Eq. (18)), as well as the internally generated variability described by  $p({\bf r}|{\bf u}^i)$  (Fig. 4a). Therefore, the distribution of  ${\bf r}$  in response to an external stimulus s has the form

$$p(\mathbf{r}|s) = \int p(\mathbf{r}|\mathbf{u}^{\mathsf{f}})p(\mathbf{u}^{\mathsf{f}}|s)d\mathbf{u}^{\mathsf{f}}.$$

For simplicity, we only compute the covariability of  $p(\mathbf{r}|\mathbf{u}^i)$  along the stimulus subspace (Fig. 1b, x-axis), because the covariability along other directions is not related with stimulus sampling. By approximating the Poissonian spiking variability  $p(\mathbf{r}|\boldsymbol{\lambda})$  with a multivariate normal distribution (Eq. (11)), and considering the limit of weak fluctuations in  $\boldsymbol{\lambda}$  along the stimulus subspace over time,  $p(\mathbf{r}|\mathbf{u}^i)$  can be computed approximately as (see math details in Supplementary Information),

$$p(\mathbf{r}|\mathbf{u}^{f}) = \int p(\mathbf{r}|\boldsymbol{\lambda})p(\boldsymbol{\lambda}|\mathbf{u}^{f})d\boldsymbol{\lambda},$$

$$\approx \mathcal{N}\left[\mathbf{r}|\mathbf{f}(s),\operatorname{diag}(\mathbf{f}(s)) + V(\bar{s}|\boldsymbol{\mu}_{f})\mathbf{f}_{s}'\mathbf{f}_{s}'^{\top}\right], \text{ where } s = \boldsymbol{\mu}_{f}.$$
(44)

 $\mathbf{f}(s) = \langle \boldsymbol{A}_t \rangle$  denotes the temporally averaged population response. The covariance structure of the neuronal response includes two terms: diag( $\mathbf{f}(s)$ ), a diagonal matrix whose entries equal that of the vector  $\mathbf{f}(s)$  denoting the (independent) Poisson spiking variability (Eq. (23)), and  $V(\bar{s}|\mu_t)\mathbf{f}_s'\mathbf{f}_s^{r\top}$ , a term that captures the covariability due to firing rate fluctuations along the stimulus subspace (Fig. 8a), where  $\mathbf{f}_s' = d\mathbf{f}(s)/ds$  is the derivative of  $\mathbf{f}(s)$  over the stimulus feature s. The covariance  $\mathbf{f}_s'\mathbf{f}_s^{r\top}$  is often termed differential (noise) correlations<sup>4,17</sup>. With the Gaussian profile of  $\mathbf{f}(s)$  (Eqs. (18) and (29)),  $\mathbf{f}_s'\mathbf{f}_s^{r\top}$  exhibits anti-symmetric structure (Fig. 8b)<sup>17,22,53,71,72</sup>.

In Eq. (44),  $V(\bar{s}|\mu_f)$  is the variance of  $\bar{s}_t$  (the mean of conditional distribution in Eq. (4a)) over time and characterizes the amplitude of internally generated differential correlations. In network implementation,  $\bar{s}_t$  and  $\mu_f$  are represented as the position of  $\lambda_t$  and  $\mathbf{u}^f$  on the

stimulus subspace respectively (Eqs. (14) and (20)). The dynamics of Gibbs sampling (Eq. S20 in Supplementary Information) and the network structure (Eq. (6)) imply that

$$V(\bar{s}|\mu_{\rm f}) = \frac{\Lambda_{\rm s}}{\Lambda_{\rm f}(\Lambda_{\rm f} + \Lambda_{\rm s})} = a^2 n_{\rm f}^{-1} w_E^*. \tag{45}$$

Note that  $V(\bar{s}|\mu_f)$  is constrained by network connections, in that it is internally generated and shared within the network (for  $w_F^* > 0$ ).

An expression for  $p(\mathbf{r}|\mathbf{s})$  can be derived similarly, and includes an additional term contributing to differential correlations compared with  $p(\mathbf{r}|\mathbf{u}^{\dagger})$  (Eq. (44)) due to fluctuations in the feedforward inputs,

$$p(\mathbf{r}|s) \approx \mathcal{N}\left[\mathbf{r}|\mathbf{f}(s), \operatorname{diag}(\mathbf{f}(s)) + V(\bar{s}|s)\mathbf{f}_{s}^{\prime}\mathbf{f}_{s}^{\prime\top}\right],$$

$$V(\bar{s}|s) = V(\bar{s}|\mu_{f}) + V(\mu_{f}|s) = \frac{\Lambda_{s}}{\Lambda_{f}(\Lambda_{f} + \Lambda_{s})} + \frac{1}{\Lambda_{f}} = a^{2}n_{f}^{-1}(w_{E}^{*} + 1).$$
(46)

Here the variance,  $V(\bar{s}|s)$ , in the stimulus feature subspace is a mixture of internal variability,  $V(\bar{s}|\mu_{\rm f})$ , and sensory noise,  $V(\mu_{\rm f}|s)$  (Eq. (23)). The neuronal response distribution in coupled networks (Fig. 6a) can be obtained similarly (see the Supplementary Information).

# A spiking network model with excitatory and inhibitory Poisson neurons

To test the proposed inference mechanisms in a network consisting of E neurons (Eqs. (10)–(37)), we simulated a well studied recurrently coupled cortical model<sup>21,22</sup>. The network consisted of  $N_E$  excitatory (E) and  $N_I$  inhibitory (I) spiking neurons, with the activity of each neuron modeled as a Hawkes process<sup>69</sup>. At time t, we represent the response of neuron j in population  $a = \{E, I\}$ ,  $\mathbf{r}_{ij}^a$ , as a spike count drawn from a Poisson distribution with instantaneous firing rate,  $\lambda_{ii}^a$ ,

$$\mathbf{r}_{tj}^{a} \sim \text{Poisson}\left[\boldsymbol{\lambda}_{tj}^{a}\right].$$
 (47)

Each neuron has a refractory period of 2ms after emitting a spike. The firing rate  $\lambda_{ij}^a$  is the sum of feedforward input  $\mathbf{u}_{ij}^{af}$  and recurrent input  $\mathbf{u}_{ij}^{af}$ , so that  $\lambda_{ij}^a = \mathbf{u}_{ij}^{af} + \mathbf{u}_{ij}^{ar}$ . The feedforward inputs are filtered spikes from upstream neurons,  $\mathbf{u}_{ij}^{af} = \sum_n \eta \left(t - t_{jn}^t\right)$ , where  $t_{jn}^t$  is the time of the nth spike received by neuron j of population a from the feedforward inputs. Here  $\eta(t)$  is the synaptic input profile which is modeled as  $\eta(t) = \exp(-t/\tau_d)/\tau_d$ , (t>0). Throughout, we set the synaptic time constant  $\tau_d = 2$ ms. To mimic the Poisson-like variability to sample a stimulus parameter in a hierarchical generative model (Eqs. (23) and (31)), the recurrent input received by neuron j in population a is defined by

$$\mathbf{u}_{tj}^{ar} = \bar{\mathbf{u}}_{tj}^{ar} + \sqrt{[\bar{\mathbf{u}}_{tj}^{ar}]_{+}} \xi_{t},$$

$$\bar{\mathbf{u}}_{tj}^{ar} = \sum_{b=1}^{N_{b}} \sum_{l=1}^{N_{b}} \frac{J_{jk}^{ab}}{\sqrt{N}} \sum_{n} \eta(t - t_{kn}^{b}),$$
(48)

where  $\bar{\mathbf{u}}_{tj}^{ar}$  is the mean recurrent input at time t given the neuronal activities of the presynaptic neurons. The recurrent input in the network is corrupted by noise whose variance equals the mean of the recurrent input. In a physiological network, recurrent noise may be generated by the chaotic state in network dynamics<sup>36</sup> or synaptic noise<sup>37,70</sup>. In Eq. (48), the function  $[\cdot]_+$  rectifies the negative input, and  $\xi_t$  is a random variable following a standard Gaussian distribution. The coefficient  $J_{ij}^{ab}$  is the synaptic weight from neuron j in population b to neuron i in population a. The time  $t_{kn}^{b}$  is the time of the nth spike fired by neuron k in population b. The parameter  $N = N_E + N_I$  is the total number of neurons in the network. The scaling of the synaptic weights by  $1/\sqrt{N}$  is standard in networks where excitation is balanced by recurrent inhibition<sup>36</sup>. Finally, the synaptic input profile of the

recurrent input,  $\eta(t)$ , is the same as the one we chose for the feedforward input for convenience. Note that the rectification in Eq. (48) on recurrent inputs will introduce errors resulting in deviations of the sampling distribution from the true posterior, and hence we chose the recurrent weights to be small (Fig. 5). The rectification only arises when using (continuous) recurrent inputs to sample the stimulus parameter, and doesn't impact the generality of sampling by (discrete) Poisson spiking variability.

To model the coding of a circular stimulus such as orientation, the excitatory neurons are arranged on a ring<sup>22,71</sup>. The preferred stimuli,  $\theta_j$ , of the excitatory neurons are equally spaced on the interval ( $-180^\circ$ ,  $180^\circ$ ], consistent with the range of latent features (Eq. (21)). Inhibitory neurons are not tuned to stimulus, and their role is to stabilize network responses. Note that the recurrent connections between E neurons are modeled using a Gaussian function decaying with the distance between the stimuli preferred by the two cells, rather than only self-connection in the simple network with only E neurons (Eq. (30)),

$$J_{jk}^{EE} = \frac{w_{EE}L}{\sqrt{2\pi}a} \exp\left[-\frac{(\theta_j - \theta_k)^2}{2a^2}\right],$$
 (49)

We imposed periodic boundaries on the Gaussian function to avoid boundary effect in simulations. Although in the generative model we assumed non-periodic feature variables (Eq. (3)), as long as the variance of the associated distributions are smaller than the width of the feature space, the network model with periodic boundaries on the recurrent connection (Eq. (49)) provides a good approximation of the non-periodic Gaussian posterior (Eq. (24)). The weight  $w_{FF}$  denotes the average connection strength of all E to E connections. The parameter  $a = 40^{\circ}$  defines the footprint of connectivity in feature space (i.e the ring), and  $L = 360^{\circ}$  is the length of the ring manifold (Eq. (21)); Multiplication by L in Eq. (49) sets the sum of all E to E connection strengths equal to  $N_E w_{EE}$ . Moreover, the excitatory and inhibitory neurons are all-to-all connected with each other (similar for I to I connections). For simplicity, we consider the E to I, I to I and I to E connections all to be unstructured (in feature space) and assume that connections of the same type have equal weight, i.e.,  $J_{jk}^{EI} = w_{EI}$ ,  $J_{jk}^{IE} = w_{EI}$ , and  $J_{jk}^{II} = w_{II}$ . To simplify the network further, we consider the connections from the same population of neurons to have the same average weight, i.e.,  $w_{EE} = w_{IE} \equiv w_E$  and  $w_{II} = w_{EI} \equiv w_I$ . For the feedforward network model shown in Fig. 2, we only remove the E recurrent connections between E neurons, i.e.,  $w_{EE} = 0$ , while keeping other connections, including  $w_{EI}$ ,  $w_{II}$ , and  $w_{IE}$ , the same as the recurrent network.

The feedforward inputs applied to E neurons consist of independent Poisson spike counts as described by Eq. (18), with rate  $\langle \mathbf{u}_f^F(s) \rangle = U^f e^{-(s-\theta_J)^2/(4a^2)}$ . The inhibitory neurons also receive feedforward indpendent Poissonian inputs. The firing rate of the input received by every *I* neuron is proportional to the overall feedforward rate of input to *E* neurons, in order to keep the excitatory and inhibitory balance of neuronal activities in the network,

$$\langle \mathbf{u}_{j}^{f} \rangle = \frac{w_{ff}}{N_{I}} \sum_{i=1}^{N_{E}} \langle \mathbf{u}_{j}^{Ef}(s) \rangle. \tag{50}$$

In the simulations, we started with a network of  $N_E$ =180 excitatory and  $N_I$ =45 inhibitory neurons, and increased the number of neurons by a fixed factor in Fig. 1d. The ratio between the average connection from I neurons and the one from E neurons was kept fixed with  $w_I/w_E$ =5. We set the feedforward weight of input to I neurons to  $w_I$ =0.8. We simulated the dynamics of the model network using the Euler method with a time step of 0.1ms. The typical parameters used in simulation can be found in Table 1 in Supplementary Information. Further details about the simulations and numerical estimates of

mutual information and linear Fisher information are also presented in Supplementary Information. The code of network simulation was written in MATLAB 2018b, and can be found at GitHub (https://github.com/wenhao-z/Sampling PoissSpk Neuron).

A spiking network model of coupled neural circuits. In the coupled neural circuits used to infer latent variables organized in parallel (Fig. 6a) the two networks are copies of each other, i.e., the two networks have the same intrinsic parameters. Each network is equivalent to the one described in the previous section, except that there is no recurrent connections between E neurons in the same network, and no variability in recurrent interactions (no noise in Eq. (48)). The absence of recurrent connections between E neurons in the same network is due to the uniform marginal prior of stimulus. Nevertheless, in the same network the E and I neurons are connected using the same connection profile as above to keep network activity stable. Between the two networks, there are only E connections which target both E and I neurons. The connections between E neurons across networks have the same pattern as that given described by Eq. (49) with the peak connection strength from network n to network m denoted as  $w_{FF}^{mn}$ . The connections from E neurons in one network to I neurons in the other is set to the same as the peak strength of E connections across networks for simplicity, i.e.,  $w_{IE}^{mn} = w_{EE}^{mn}$ . To simplify the network model further, we set the inter-network connections to be symmetric, which means  $w_{FF}^{mn} = w_{FF}^{nm}$ . In the simulations  $w_{FF}^{mn}$  was adjusted to determine how the sampling distribution is affected (Fig. 7a).

Comparing the sampling distribution with posterior in coupled neural circuits. We read out the samples from the posterior distribution of each stimulus,  $\tilde{s}_{mt}$ , individually from the spiking activities of E neurons,  $\mathbf{r}_{mt}$ , in network m in every time window of 20ms by using a population vector. We used this collection of samples to estimate the mean,  $\langle \tilde{\mathbf{s}} \rangle = (\langle \tilde{s}_1 \rangle, \langle \tilde{s}_2 \rangle)^{\top}$ , and covariance matrix,  $\Sigma_{\mathbf{s}}$ , of the sampling distribution. Meanwhile, the mean  $\mu_{\rm f}$  and precision matrix  $\Lambda_{\rm f}$  of the likelihood are linearly read out from the feedforward inputs fed into the network model (Eq. (33)).

If the sampling distribution is comparable with the posterior, the sampling mean  $\langle \tilde{\mathbf{s}} \rangle$  and covariance  $\mathbf{\Sigma}_s$  should satisfy Eq. (34). We use the actual sampling covariance and the likelihood parameters to predict the sampling mean, i.e.,  $\langle \tilde{\mathbf{s}} \rangle_{pred} = \mathbf{\Sigma}_s \mathbf{\Lambda}_f \boldsymbol{\mu}_f$ , and compare it with the actual  $\langle \tilde{\mathbf{s}} \rangle$  (Fig. 7d–f). To obtain the posterior precision matrix, given the sampling mean  $\langle \tilde{\mathbf{s}} \rangle$  and the likelihood parameters, we vary the single parameter of prior precision  $\Lambda_s$  to minimize the KL divergence from the prediction of posterior by using the value of  $\Lambda_s$ , and the actual sampling distribution. Given this value of  $\Lambda_s$ , the prediction of posterior precision is computed as  $\mathbf{K}_{pred} = \mathbf{\Lambda}_s + \mathbf{\Lambda}_f$  (Eq. (34)) which is then compared with actual sampling precision matrix ( $\mathbf{\Sigma}_s^{-1}$ ; see Fig. 7c–g). The prior precision,  $\Lambda_s$ , is a *subjective* prior, which reflects the prior stored in the recurrent network and may change with input (see Discussion). More details of network simulation and parameters can be found in Supplementary Information.

# Reporting summary

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

# **Data availability**

This is a strictly computational study and all data used in making figures were generated by computer simulations of the proposed model with the link of codes shown in Code Availability.

#### Code availability

The code of network simulation was written in MATLAB 2018b, and can be found at GitHub (https://github.com/wenhao-z/Sampling\_PoissSpk\_Neuron)<sup>73</sup>.

### References

- Pouget, A., Dayan, P. & Zemel, R. S. Inference and computation with population codes. *Ann. Rev. Neurosci.* 26, 381–410 (2003).
- Doiron, B., Litwin-Kumar, A., Rosenbaum, R., Ocker, G. K. & Josić, K. The mechanics of state-dependent neural correlations. *Nat. Neurosci.* 19, 383–393 (2016).
- 3. Goris, R. L., Movshon, J. A. & Simoncelli, E. P. Partitioning neuronal variability. *Nat. Neurosci.* **17**, 858–865 (2014).
- Kohn, A., Coen-Cagli, R., Kanitscheider, I. & Pouget, A. Correlations and neuronal population information. *Ann. Rev. Neurosci.* 39, 237–256 (2016).
- Harris, J. A. et al. Hierarchical organization of cortical and thalamic connectivity. *Nature* 575, 195–202 (2019).
- Oh, S. W. et al. A mesoscale connectome of the mouse brain. Nature 508, 207–214 (2014).
- Douglas, R. J. & Martin, K. A. Neuronal circuits of the neocortex. Annu. Rev. Neurosci. 27, 419–451 (2004).
- Rossi, L. F., Harris, K. D. & Carandini, M. Spatial connectivity matches direction selectivity in visual cortex. *Nature* 588, 648–652 (2020).
- Harris, K. D. & Mrsic-Flogel, T. D. Cortical connectivity and sensory coding. Nature 503, 51–58 (2013).
- Ernst, M. O. & Banks, M. S. Humans integrate visual and haptic information in a statistically optimal fashion. *Nature* 415, 429–433 (2002).
- Yuille, A. & Kersten, D. Vision as Bayesian inference: analysis by synthesis? *Trends Cogn. Sci.* 10, 301–308 (2006).
- Lee, T. S. & Mumford, D. Hierarchical Bayesian inference in the visual cortex. JOSA A 20, 1434–1448 (2003).
- Körding, K. P. & Wolpert, D. M. Bayesian integration in sensorimotor learning. *Nature* 427, 244–247 (2004).
- Knill, D. C. & Pouget, A. The Bayesian brain: the role of uncertainty in neural coding and computation. *Trends Neurosci.* 27, 712–719 (2004).
- Pouget, A., Beck, J. M., Ma, W. J. & Latham, P. E. Probabilistic brains: knowns and unknowns. *Nat. Neurosci.* 16, 1170 (2013).
- Fiser, J., Berkes, P., Orbán, G. & Lengyel, M. Statistically optimal perception and learning: from behavior to neural representations. *Trends. Cogn. Sci.* 14, 119–130 (2010).
- Moreno-Bote, R. et al. Information-limiting correlations. *Nat. Neurosci.* 17, 1410 (2014).
- Dayan, P. & Abbott, L. F. Theoretical Neuroscience, Vol. 806 (MIT Press, 2001).
- Averbeck, B. B., Latham, P. E. & Pouget, A. Neural correlations, population coding and computation. *Nat. Rev. Neurosci.* 7, 358 (2006).
- Georgopoulos, A. P., Schwartz, A. B. & Kettner, R. E. Neuronal population coding of movement direction. *Science* 233, 1416–1419 (1986).
- Rubin, D. B., Van Hooser, S. D. & Miller, K. D. The stabilized supralinear network: a unifying circuit motif underlying multi-input integration in sensory cortex. *Neuron* 85, 402–417 (2015).
- Ben-Yishai, R., Bar-Or, R. L. & Sompolinsky, H. Theory of orientation tuning in visual cortex. Proc. Natl Acad. Sci. 92, 3844–3848 (1995).
- 23. Somers, D. C., Nelson, S. B. & Sur, M. An emergent model of orientation selectivity in cat visual cortical simple cells. *J. Neurosci.* **15**, 5448–5465 (1995).
- Huang, C., Pouget, A. & Doiron, B. D. Internally generated population activity in cortical networks hinders information transmission. Sci. Adv. 8, eabg5244 (2022).
- Kersten, D., Mamassian, P. & Yuille, A. Object perception as Bayesian inference. Annu. Rev. Psychol. 55, 271–304 (2004).
- Doya, K., Ishii, S., Pouget, A. & Rao, R. P. Bayesian Brain: Probabilistic Approaches to Neural Coding (MIT press, 2007).

- Hoyer, P. O. & Hyvärinen, A. Interpreting neural response variability as Monte Carlo sampling of the posterior. In Advances in Neural Information Processing Systems, 293–300 (2003).
- Buesing, L., Bill, J., Nessler, B. & Maass, W. Neural dynamics as sampling: a model for stochastic computation in recurrent networks of spiking neurons. PLoS Comput. Biol. 7, e1002211 (2011).
- Savin, C. & Deneve, S. Spatio-temporal representations of uncertainty in spiking neural networks. In NIPS, vol. 27, 2024–2032 (2014).
- 30. Orbán, G., Berkes, P., Fiser, J. & Lengyel, M. Neural variability and sampling-based probabilistic representations in the visual cortex. *Neuron* **92**, 530–543 (2016).
- Haefner, R. M., Berkes, P. & Fiser, J. Perceptual decision-making as probabilistic inference by neural sampling. *Neuron* 90, 649–660 (2016).
- 32. Echeveste, R., Aitchison, L., Hennequin, G. & Lengyel, M. Cortical-like dynamics in recurrent circuits optimized for sampling-based probabilistic inference. *Nat. Neurosci.* **23**, 1138–1149 (2020).
- 33. Festa, D., Aschner, A., Davila, A., Kohn, A. & Coen-Cagli, R. Neuronal variability reflects probabilistic inference tuned to natural image statistics. *Nat. Commun.* **12**, 1–11 (2021).
- 34. Hénaff, O. J., Boundy-Singer, Z. M., Meding, K., Ziemba, C. M. & Goris, R. L. Representation of visual uncertainty through neural gain variability. *Nat. Commun.* 11, 1–12 (2020).
- 35. 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).
- Vreeswijk, C. V. & Sompolinsky, H. Chaotic balanced state in a model of cortical circuits. Neural Comput. 10, 1321–1371 (1998).
- Rosenbaum, R., Rubin, J. & Doiron, B. Short term synaptic depression imposes a frequency dependent filter on synaptic information transfer. PLoS Comput. Biol. 8, e1002557 (2012).
- 38. Bishop, C. M. Pattern Recognition and Machine Learning (Springer, 2006).
- Ma, W. J., Beck, J. M., Latham, P. E. & Pouget, A. Bayesian inference with probabilistic population codes. *Nat. Neurosci.* 9, 1432–1438 (2006).
- Jazayeri, M. & Movshon, J. A. Optimal representation of sensory information by neural populations. *Nat. Neurosci.* 9, 690–696 (2006).
- Lewicki, M. S. & Sejnowski, T. J. Bayesian unsupervised learning of higher order structure. Adv. Neural Inf. Process. Syst. 9, 529–535 (1996).
- 42. Grabska-Barwinska, A., Beck, J. M., Pouget, A. & Latham, P. E. Demixing odors-fast inference in olfaction. *Adv. Neural Inf. Process.* Syst. **26**, 1–9 (2013).
- 43. Field, D. J., Hayes, A. & Hess, R. F. Contour integration by the human visual system: evidence for a local "association field". *Vision Res.* **33**, 173–193 (1993).
- Geisler, W. S., Perry, J. S., Super, B. & Gallogly, D. Edge cooccurrence in natural images predicts contour grouping performance. Vision Res. 41, 711–724 (2001).
- 45. Cossell, L. et al. Functional organization of excitatory synaptic strength in primary visual cortex. *Nature* **518**, 399–403 (2015).
- Kanitscheider, I., Coen-Cagli, R., Kohn, A. & Pouget, A. Measuring fisher information accurately in correlated neural populations. *PLoS Comput. Biol.* 11, e1004218 (2015).
- 47. Lee, T. S. The visual system's internal model of the world. *Proc. IEEE* **103**, 1359–1378 (2015).
- 48. Vasudeva Raju, R. & Pitkow, Z. Inference by reparameterization in neural population codes. *Adv. Neural Inf. Process. Syst.* **29**, 2029–2037 (2016).
- Beck, J. M., Latham, P. E. & Pouget, A. Marginalization in neural circuits with divisive normalization. *J. Neurosci.* 31, 15310–15319 (2011).

- Aitchison, L. & Lengyel, M. The Hamiltonian brain: Efficient probabilistic inference with excitatory-inhibitory neural circuit dynamics. PLoS Comput. Biol. 12, e1005186 (2016).
- Shivkumar, S., Lange, R., Chattoraj, A. & Haefner, R. A probabilistic population code based on neural samples. In *Advances in Neural Information Processing Systems*, Vol. 31 (eds Bengio, S. et al.) (Curran Associates, Inc., 2018). https://proceedings.neurips.cc/paper/2018/ file/5401acfe633e6817b508b84d23686743-Paper.pdf.
- Kanitscheider, I., Coen-Cagli, R. & Pouget, A. Origin of informationlimiting noise correlations. *Proc. Natl Acad. Sci.* 112, E6973–E6982 (2015).
- Ponce-Alvarez, A., Thiele, A., Albright, T. D., Stoner, G. R. & Deco, G. Stimulus-dependent variability and noise correlations in cortical mt neurons. *Proc. Natl Acad. Sci.* 110, 13162–13167 (2013).
- Wu, S., Wong, K. M., Fung, C. A., Mi, Y. & Zhang, W. Continuous attractor neural networks: candidate of a canonical model for neural information representation. *F1000Research* 5, F1000 (2016).
- Hennequin, G., Ahmadian, Y., Rubin, D. B., Lengyel, M. & Miller, K. D. The dynamical regime of sensory cortex: stable dynamics around a single stimulus-tuned attractor account for patterns of noise variability. *Neuron* 98, 846–860 (2018).
- Lange, R. D. & Haefner, R. M. Task-induced neural covariability as a signature of approximate Bayesian learning and inference. PLoS Comput Biol 18, e1009557 (2022).
- 57. Zhang, K. Representation of spatial orientation by the intrinsic dynamics of the head-direction cell ensemble: a theory. *J. Neurosci.* **16**, 2112–2126 (1996).
- Deneve, S., Latham, P. E. & Pouget, A. Reading population codes: a neural implementation of ideal observers. *Nat. Neurosci.* 2, 740–745 (1999).
- Wu, S., Amari, S.-i & Nakahara, H. Population coding and decoding in a neural field: a computational study. *Neural Comput.* 14, 999–1026 (2002).
- Schulz, E., Tenenbaum, J. B., Duvenaud, D., Speekenbrink, M. & Gershman, S. J. Compositional inductive biases in function learning. Cogn. Psychol. 99, 44–79 (2017).
- Abbott, L. F., Varela, J., Sen, K. & Nelson, S. Synaptic depression and cortical gain control. Science 275, 221–224 (1997).
- 62. Ermentrout, B. Linearization of fi curves by adaptation. *Neural Comput.* **10**, 1721–1729 (1998).
- 63. Coen-Cagli, R., Kohn, A. & Schwartz, O. Flexible gating of contextual influences in natural vision. *Nat. Neurosci.* **18**, 1648–1655 (2015).
- Beck, J. M., Ma, W. J., Pitkow, X., Latham, P. E. & Pouget, A. Not noisy, just wrong: the role of suboptimal inference in behavioral variability. Neuron 74, 30–39 (2012).
- 65. Churchland, M. M. et al. Stimulus onset quenches neural variability: a widespread cortical phenomenon. *Nat. Neurosci.* **13**, 369–378 (2010).
- Maimon, G. & Assad, J. A. Beyond Poisson: increased spike-time regularity across primate parietal cortex. *Neuron* 62, 426–440 (2009).
- 67. Zhang, W., Lee, T. S., Doiron, B. & Wu, S. Distributed sampling-based Bayesian inference in coupled neural circuits. *bioRxiv* (2020).
- 68. Ganguli, D. & Simoncelli, E. P. Implicit encoding of prior probabilities in optimal neural populations. *Adv. Neural Inf. Process. Syst.* **2010**, 658 (2010).
- Trousdale, J., Hu, Y., Shea-Brown, E. & Josić, K. Impact of network structure and cellular response on spike time correlations. *PLoS Comput. Biol.* 8, e1002408 (2012).

- Rusakov, D. A., Savtchenko, L. P. & Latham, P. E. Noisy synaptic conductance: Bug or a feature? *Trends Neurosci.* 43, 363–372 (2020).
- 71. Wu, S., Hamaguchi, K. & Amari, S.-i Dynamics and computation of continuous attractors. *Neural Comput.* **20**, 994–1025 (2008).
- Wimmer, K., Nykamp, D. Q., Constantinidis, C. & Compte, A. Bump attractor dynamics in prefrontal cortex explains behavioral precision in spatial working memory. *Nat. Neurosci.* 17, 431 (2014).
- 73. Zhang, W.-H. Sampling-based Bayesian inference in recurrent circuits of stochastic spiking neurons. Sampling\_PoissSpk\_Neuron. https://doi.org/10.5281/zenodo.8088755 (2023).

# Acknowledgements

National Institutes of Health grants 1R01MH115557 (K.J.), 1U19NS107613-01 (B.D.), R01EB026953 (B.D.), R01EY034723 (B.D.); National Science Foundation grant NSF-DBI-1707400 (K.J.), DMS-2207647 (K.J.). Vannevar Bush faculty fellowship N00014-18-1-2002 (B.D); Simons Foundation Collaboration on the Global Brain (B.D.).

## **Author contributions**

W.H.Z., S.W., K.J., and B.D. conceived and designed the study. W.H.Z. developed the ideas, performed the analyses, and the numerical simulations. W.H.Z, K.J., and B.D. discussed the results and wrote the manuscript.

## **Competing interests**

The authors declare no competing interests.

## **Additional information**

**Supplementary information** The online version contains supplementary material available at https://doi.org/10.1038/s41467-023-41743-3.

**Correspondence** and requests for materials should be addressed to Krešimir Josić or Brent Doiron.

**Peer review information** *Nature Communications* thanks the anonymous reviewer(s) for their contribution to the peer review of this work. A peer review file is available.

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

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

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit http://creativecommons.org/licenses/by/4.0/.

© The Author(s) 2023