

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

Wen-Hao Zhang^{1-4,†}, Si Wu⁵, Krešimir Josić^{6,7*} Brent Doiron^{1-4*}

¹Departments of Neurobiology and Statistics, University of Chicago, Chicago, IL, USA.

²Grossman Center for Quantitative Biology and Human Behavior, University of Chicago, Chicago, IL, USA.

³Department of Mathematics, University of Pittsburgh, Pittsburgh, PA, USA.

⁴Center for the Neural Basis of Cognition, Pittsburgh, PA, USA.

⁵School of Electronics Engineering and Computer Science, IDG/McGovern Institute for Brain Research, Peking-Tsinghua Center for Life Sciences, Peking University, Beijing 100871, China.

⁶Department of Mathematics, University of Houston, TX, USA.

⁷Department of Biology and Biochemistry, University of Houston, TX, USA.

† Present address: Lyda Hill Department of Bioinformatics,
UT Southwestern Medical Center, Dallas, TX, USA.

*Equal contribution.

Correspondence should be addressed to: bdoiron@uchicago.edu or kresimir.josic@gmail.com

Abstract

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.

Keywords: Sampling-based Bayesian inference, Poisson spiking neurons, Recurrent network dynamics, Differential correlations.

16 Introduction

17 In an uncertain and changing world, it is imperative for the brain to reliably represent and interpret
18 external stimuli. The cortex is essential for the representation of the sensory world, and it is believed
19 that populations of neurons collectively code for richly structured sensory scenes [1]. However,
20 two central characteristics of cortical circuits remain to be properly integrated into population
21 coding frameworks. First, neuronal activity in sensory cortices is often noisy, showing significant
22 variability of spiking responses evoked by the same stimulus [2, 3]. In many traditional coding
23 frameworks such spiking variability degrades the representation of stimuli by cortical activity [4].
24 Why cortical responses display large spiking variability while isolated cortical neurons can respond
25 reliably remains a mystery. Second, the primary source of synaptic inputs to cortical neurons
26 does not come from upstream centers which convey sensory signals, but rather from recurrent
27 pathways between cortical neurons [5–7]. While such recurrent connections are often organized
28 about a stimulus feature axis [8, 9], it is not obvious whether or how their presence improves
29 overall representation. We propose a biologically motivated inference coding scheme where these
30 two ubiquitous cortical circuit features, variability in spike generation and recurrent connections,
31 together support a probabilistic representation of stimuli in rich sensory scenes.

32 Numerous studies have framed sensory processing in the cortex in terms of Bayesian inference
33 (e.g., [10–16]). Specifically, the ‘Bayesian brain’ hypothesis posits that sensory cortex infers and
34 synthesizes a posterior distribution of the latent stimuli which describes the probability of possible
35 stimuli that could have given rise to the sensory inputs. Performing Bayesian inference requires cor-
36 tex to store an internal model that represents how sensory inputs and external stimuli are generated.
37 Once a sensory input is received, cortical dynamics inverts this internal model in a process termed
38 ‘analysis-by-synthesis’ [12], and represents the posterior distributively across neurons and/or across
39 time [15, 16]. In this study, we propose that recurrent connections in cortical circuits store the prior
40 of latent stimuli to produce the posterior distribution when combined with evidence from sensory
41 inputs. Moreover, we posit that Poisson spiking variability provides a source of fluctuations needed
42 for generating random samples from the inferred posterior.

43 To test these hypotheses we consider a recurrent circuit model where neurons receive stochastic
44 feedforward inputs which carry information about the external world, and respond with Poisson-
45 distributed spiking activity. We find that such Poissonian spiking provides the variability that allows
46 the network to generate samples from posterior stimulus distributions with differing uncertainties.
47 We use this sampling framework to illustrate circuit-based Bayesian inference given two distinct
48 generative models of stimuli in the external world: one organized hierarchically with a stimulus
49 variable that depends on a latent context variable, and a second where a pair of latent stimuli are
50 organized in parallel. In both cases a recurrent circuit is able to generate samples from the joint
51 posterior, and infer the values of the latent variables. We show through both analytic derivation

52 and simulations that recurrent connections represent the correlation structure of these models, and
53 the weight of these connections can be tuned to optimally capture the prior distribution of stimuli
54 in the external world. The stronger the correlation between the latent variables, the stronger the
55 recurrent connections need to be for the network to generate samples from the correct posterior
56 distribution.

57 Finally, a neural signature of this circuit-based sampling mechanism is internally generated
58 population noise correlations aligned with the stimulus response direction, often referred to as “dif-
59 ferential correlations” [4, 17]. In our framework, the amplitude of internally generated differential
60 correlations is determined by the recurrent connection strength, which also determines the prior
61 stored by the circuit. Since optimal inference requires a specific magnitude of recurrent connectiv-
62 ity, differential correlations resulting from such recurrent connectivity are a potential signature of
63 optimal coding. This is in contrast to the deleterious impact of externally generated differential
64 correlations. We thus predict that the correlation structure of the external world shapes recurrent
65 wiring in neural circuits, and is reflected in the pattern of differential noise correlations. We use
66 this logic to provide testable predictions from our framework for sampling-based Bayesian inference
67 by recurrent, stochastic cortical circuits.

68 Results

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

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

76 In population coding frameworks stimuli are encoded by a neuronal population with individual
77 neurons tuned to a preferred stimulus value. The preferred values of all neurons cover the whole
78 range of stimuli [18–20] (Fig. 1B, bottom); if s ranges over a periodic domain (such as the orientation
79 of a bar in a visual scene, or the direction of an arm reach) then it is commonly assumed that the
80 neurons’ preferred stimuli are distributed on a ring (Fig. 1B, top). To generate neuronal responses
81 from such a population we simulate a network of neurons whose spiking activity, \mathbf{r}_t , at time t is
82 Poissonian with instantaneous firing rate λ_t (Eq. 11). For simplicity we assume linear (or linearized)
83 neuronal transfer and synaptic interactions (Eqs. 10–11), so that the firing rate is a linear function
84 of the feedforward and recurrent inputs. We couple excitatory (E) neurons with similar stimulus
85 preferences more strongly [8, 9] to one another, compared to neuron pairs with dissimilar tuning. In

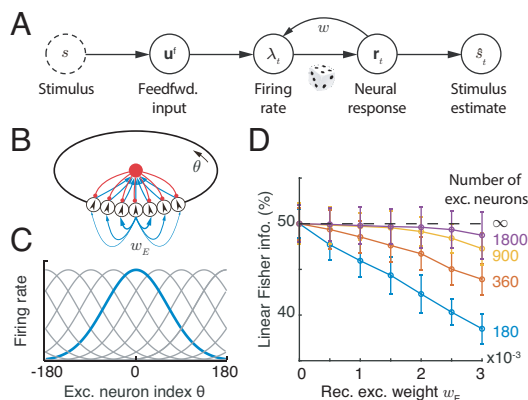


Figure 1: A network with structured recurrent connections limits the linear Fisher Information (LFI) about external stimuli. (A) A schematic diagram showing how a stimulus, s , is encoded in neuronal response, \mathbf{r}_t . A stimulus estimate, \hat{s}_t , can be obtained from \mathbf{r}_t . (B) A recurrent 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.

86 this way the recurrent E connectivity has the same circular symmetry as the stimulus (Fig. 1B). In
 87 contrast, connections between inhibitory (I) neurons are unstructured, and inhibitory activity acts
 88 to stabilize network activity [21]. A stimulus, e.g. $s = 0$, results in elevated activity of E neurons
 89 with the corresponding preference (Fig. S1A). As expected, an increase in the strength of recurrent
 90 excitatory connections increases both the firing rates and the trial-to-trial pairwise covariability
 91 (i.e. noise correlations) in the responses [2] (Fig. S2A). This canonical network model has been
 92 widely used to explain cortical network dynamics and neural coding [21–23].

93 We use linear Fisher Information (LFI) to quantify the impact of recurrent connectivity and
 94 internal spiking variability on the accuracy of the stimulus estimate, \hat{s}_t , from the activity vector \mathbf{r}_t
 95 (see details in Eq. S39 in Supplemental Information). The inverse of LFI provides a lower bound
 96 on the expected square of the difference between the true value, s , and the estimate, \hat{s}_t , made by a
 97 linear decoder [1, 4, 17–19, 24]. In the limit of an infinite number of neurons available to the decoder
 98 LFI is unaffected by recurrent connectivity strength, w_E (Fig. 1D, dashed line). This is because
 99 the mean response of the network is linear in its inputs, and an (invertible) linear transformation
 100 can neither increase nor decrease LFI (see Eq. S38 in Supplemental Information). For networks
 101 with a finite number of neurons, the variability from spike generation is shared between neurons
 102 via recurrent interactions. Consequently an increase in coupling strength, w_E , reduces LFI in finite
 103 networks (Fig. 1D, colored lines).

104 In sum, recurrent connectivity and spiking variability do not improve, and often degrade, stim-

105 ulus representation in the network (as measured by LFI). Since synaptic coupling is biologically
106 expensive, a network that most accurately and cheaply represents a stimulus is then one with no
107 recurrent connections (i.e., $w_E = 0$) and minimal spiking variability. Nevertheless, connectivity
108 in mammalian cortex is highly recurrent [5–7, 9], and neural responses are highly variable [2, 3].
109 What is then the purpose of these extensive recurrent connections between cortical neurons, and
110 why are their responses so noisy?

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

119 Internally generated Poisson spiking variability drives sampling-based Bayesian infer- 120 ence

121 Many studies suggest that neuronal response variability is a signature of sampling in neural circuits
122 (e.g., [16, 27–32]). In these studies the instantaneous population responses, \mathbf{r}_t , represent a sample
123 of a latent stimulus, and the empirical distribution of stimulus samples collected over time is an
124 approximation of the posterior distribution. Furthermore, response variability is typically modeled
125 using a continuous (e.g., Gaussian) distribution [27, 29–33]. However, spike trains from cortical
126 neurons are often Poissonian, and spike counts are discrete [3, 34]. It is unclear if discrete Poisso-
127 nian variability can generate samples from stimuli with continuous probability distributions (e.g.,
128 orientation, moving direction) with the flexibility needed to represent different stimulus uncertain-
129 ties.

130 We address this question using a theory based on a simple model network composed of excitatory
131 (E) Poissonian neurons (Eqs. 10–11), and subsequently support our findings by simulating a network
132 containing both E and inhibitory (I) neurons (e.g. Fig. 1B). We start by showing that Poissonian
133 spiking in a population of tuned neurons can drive sampling from a well-defined distribution.
134 We assume that the instantaneous firing rates of a population of E neurons, $\boldsymbol{\lambda}_t$, have a bell-shaped
135 (Gaussian) profile (Fig. 2B), so that for the j^{th} neuron $\lambda_{tj} = R \exp[\mathbf{h}_j(\bar{s}_t)] = R \exp[-(\bar{s}_t - \theta_j)^2/2a^2]$
136 (See Eq. 12 in Methods). Here θ_j is the preferred stimulus of neuron j , a is the width of the tuning
137 curve, and \bar{s}_t is the location of the peak of the firing rate profile, $\boldsymbol{\lambda}_t$, in stimulus space (x-axis in
138 Fig. 2B). Note that the value of \bar{s}_t is arbitrary here, but we will later relate it to the input to the
139 population. The (smooth) Gaussian tuning curves simplify the analysis, but are not essential for

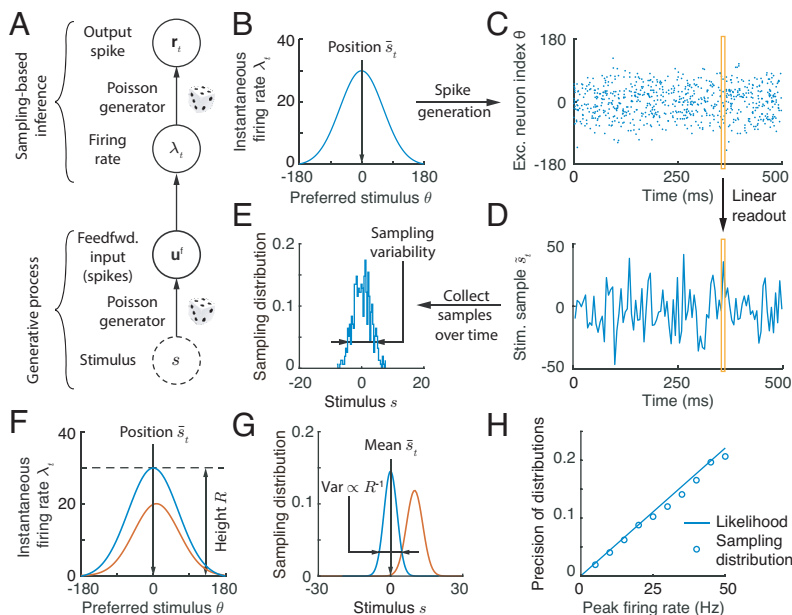


Figure 2: Spike generation with Poissonian variability can support sampling-based Bayesian inference. (A) We use a feedforward network model (no recurrent connections) to demonstrate how spiking variability drives sampling. Neurons receive feedforward inputs, u^f , modeled as independent Poisson spike trains, resulting in a Poissonian population response, r_t , with means determined by the instantaneous firing rate vector, λ_t . (B-E) Demonstration of sampling via stochastic spike generation. A population of neurons with Gaussian tuning and firing rates λ_t (B) generates a realization a population response, r_t (C). A sample from the posterior distribution of the stimulus (D, orange box) can be linearly read out from the population response (C, orange box). (E) The sampling distribution is obtained by collecting stimulus samples over time. (F-G) The profile of population firing rates (F) determines the sampling distribution (G). The position of the population firing rate, \bar{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. (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.

140 the argument. Finally, the preferred stimuli of the E neurons, $\{\theta_j\}_{j=1}^{N_E}$, are uniformly distributed
 141 over the stimulus range (Fig 1B). In each time interval the population activity is given by a vector
 142 of independent Poisson random variables, r_t , with means determined by the instantaneous firing
 143 rate vector λ_t (Fig. 2B-C). At each time, t , this spiking activity produces a stimulus sample, \tilde{s}_t ,
 144 from the probability distribution determined by the instantaneous firing rates, λ_t (Fig. 2D, see
 145 Methods),

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

146 With the Gaussian firing rate profile we use here, the stimulus sample, \tilde{s}_t , can be read out as
 147 $\tilde{s}_t = \sum_j r_{tj} \theta_j / \sum_j r_{tj}$ (Eq. 14 and Fig. 2D), which can be thought of as the location of the response,
 148 r_t , in stimulus space (y-axis in Fig. 2C). The collection of stimulus samples across time ($\{\tilde{s}_t\}$;
 149 Fig. 2E), determines the sampling distribution $q(s) = T^{-1} \sum_t \delta(s - \tilde{s}_t)$ which approximates the

150 distribution $p(s|\boldsymbol{\lambda}_t)$, i.e., $p(s|\boldsymbol{\lambda}_t) \approx q(s)$ [16, 35]. Here $\delta(\cdot)$ is the Dirac delta function and T is the
 151 number of samples.

152 To use this mechanism to produce samples from the posterior distribution of a stimulus, we
 153 must define a generative model for the feedforward inputs evoked by a stimulus. We take the
 154 feedforward input to the neural population, \mathbf{u}^f , to be a vector of independent Poisson spike counts
 155 with Gaussian tuning over the stimulus, s . Following assumptions widely used in previous studies
 156 of probabilistic population codes (PPC) [36, 37], we assume that the mean input spike count to
 157 the j^{th} excitatory neuron in the population is $\langle \mathbf{u}_j^f(s) \rangle \propto \exp[\mathbf{h}_j(s)] = \exp[-(s - \theta_j)^2/2a^2]$. A
 158 single realization of the input, \mathbf{u}^f , in a time interval encodes the whole likelihood function over the
 159 stimulus, $p(\mathbf{u}^f|s)$ [36]. This likelihood is proportional to a Gaussian due to the Gaussian profile of
 160 feedforward input (Eq. 19),

$$\begin{aligned} p(\mathbf{u}^f|s) &= \prod_{j=1}^{N_E} \text{Poisson}[\langle \mathbf{u}_j^f(s) \rangle], \\ &\propto \exp[\mathbf{h}(s)^\top \mathbf{u}^f], \\ &\propto \mathcal{N}(s|\mu_f, \Lambda_f^{-1}). \end{aligned} \tag{2}$$

161 Here the likelihood mean, μ_f , is determined by the location of \mathbf{u}^f in stimulus space, and the
 162 precision, Λ_f , is proportional to the spike count (or height) of \mathbf{u}^f (Eq. 20). Since a realization of
 163 the feedforward input encodes the whole likelihood function, we present a fixed \mathbf{u}^f to the network
 164 over time (dropping the time index t), and describe how samples from the posterior $p(s|\mathbf{u}^f)$ are
 165 generated by the network.

166 A simple example of inference via sampling is provided by a population of E neurons with-
 167 out recurrent connections and instantaneous firing rates equal to the feedforward input, $\boldsymbol{\lambda}_t = \mathbf{u}^f$
 168 (Eq. 10), and hence constant in time (Fig. 2A). In this feedforward network Poisson spike generation
 169 produces samples from the normalized likelihood, i.e., $\tilde{s}_t \sim p(\tilde{s}|\boldsymbol{\lambda}_t) \propto p(\mathbf{u}^f|\tilde{s})$, and consequently the
 170 network represents a uniform stimulus prior (i.e., $p(s)$ is a constant).

171 To test our theory, we simulated the response of a network of tuned excitatory (E) and untuned
 172 inhibitory (I) neurons (Fig. 2A,C) to a fixed but randomly generated feedforward input (Eq. 18).
 173 While the E neurons shared no recurrent connections, the E and I neurons were connected to main-
 174 tain stable network activity. To confirm that the overall firing rate dictated the sampling variability
 175 (Eq. 1), we increased the feedforward input rate, which reduced the width of the likelihood (Eq. 2).
 176 As a result, the sampling precision (inverse of the sampling variance) increased and matched the
 177 precision of the likelihood (Fig. 2G, H), even as the normalized response variability (measured the
 178 by Fano factor) of single neurons remained unchanged.

179 While the above analysis introduces the key components of a sampling-based theory of inference,
 180 stimulus sampling using a feedforward network is unnecessary: A single observation of the response
 181 \mathbf{r} in a deterministic feedforward network ($\mathbf{r} = \mathbf{u}^f$ after removing spike generation in Eq. 11) would

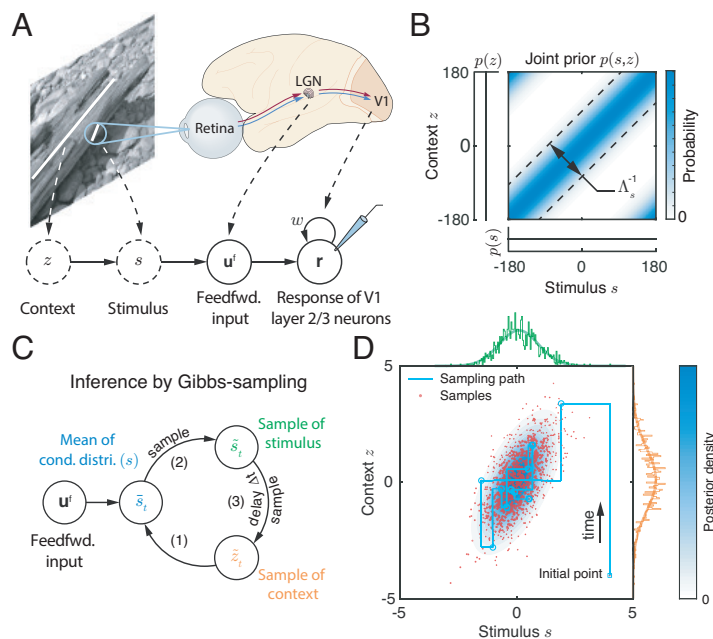


Figure 3: A hierarchical generative model and posterior inference via Gibbs sampling. (A) An example of sensory feedforward input generation: The context, z , is the orientation of the tree trunk, and the stimulus, s , is the orientation of the bark texture located in the classical receptive field of a V1 hypercolumn. The recurrent circuit generates samples from the joint posterior over stimulus and context. Solid circles: observations and responses in the brain; dashed circles: latent variables in the external world. (Natural image and brain schematic adapted from [38] and [39] respectively). (B) The joint prior over the context, z , and stimulus, s , is concentrated on the diagonal. The correlation between context and stimulus is determined by parameter Λ_s . (C) The posterior over context 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 the latent stimulus, s , and context, 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 the stimulus s and context z respectively.

182 also represent the whole likelihood [36], avoiding the costly process of collecting samples \tilde{s}_t across
 183 time. We next consider more interesting cases, and show that spiking variability in recurrent
 184 networks can drive sampling from more complex posterior distributions.

185 Recurrent cortical circuit samples a hierarchical generative model

186 Recurrent networks can store a variety of generative model structures; to demonstrate the generality
 187 of our sampling framework we provide two example generative models which serve as building blocks
 188 for more complex models. We first consider a two-stage hierarchical model of feedforward inputs
 189 received by the cortical circuit (Fig. 3A). The first stage of our model consists of a stimulus, s ,
 190 and a context, z , both of which are one dimensional for simplicity. The structure of the world
 191 is described by the joint distribution, $p(s,z)$. Using the visual system as motivation, s , could
 192 be the orientation of the visual texture within a classical receptive field (local information) of a

193 hypercolumn of V1 neurons, while the orientation within a non-classical receptive field of these
 194 cells could describe the corresponding context, z (Fig. 3A). The likelihood of the stimulus based
 195 on a given context, $p(s|z) = \mathcal{N}(s|z, \Lambda_s^{-1})$, is Gaussian with precision Λ_s . For simplicity, we assume
 196 that the context prior, $p(z)$, is uniform, which implies that the marginal prior of s , is also uniform
 197 (Fig. 3B). This assumption is not essential for our main conclusions but does simplify the analysis.
 198 Importantly, the joint prior of stimulus and context, $p(s, z)$, can have non-trivial structure with
 199 the density concentrated around the diagonal $s = z$ (Fig. 3B). The precision Λ_s measures how
 200 strongly the context, z , and the stimulus, s , are related, and thus determines how strongly their
 201 joint distribution is concentrated around the diagonal.

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

$$\begin{aligned} p(\mathbf{u}^f | s) p(s | z) p(z) &\propto \prod_{j=1}^{N_E} \text{Poisson}(\mathbf{u}_j^f | s) p(s | z), \\ &\propto \mathcal{N}(s | \mu_f, \Lambda_f^{-1}) \mathcal{N}(s | z, \Lambda_s^{-1}). \end{aligned} \quad (3)$$

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

209 Gibbs sampling of the joint stimulus and context posterior

210 One approach to approximate the joint distribution over stimulus and context is Gibbs sampling [31,
 211 35, 40, 41] which starts with an initial guess for the value of the two latent variables, and proceeds
 212 by alternately generating samples of one variable from the distribution conditioned on the value of
 213 the second variable. More precisely, to approximate the joint posterior of s and z (Eq. 3), Gibbs
 214 sampling proceeds by generating a sequence of samples, $(\tilde{s}_t, \tilde{z}_t)$ indexed by time t , through recursive
 215 iteration of the following steps (Fig. 3C and Eq. 25),

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

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

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

216 Here Δt is the time increment between successive samples. The samples (red dots in Fig. 3D) are
 217 generated by alternately fixing the values of the two variables, so that sampling trajectories alternate
 218 between horizontal and vertical jumps (cyan lines in Fig. 3D). The empirical distribution of samples,
 219 i.e., $q(s, z | \mathbf{u}^f) = T^{-1} \sum_t \delta[(s, z)^\top - (\tilde{s}_t, \tilde{z}_t)^\top]$ with \top denoting vector transpose, approximates the

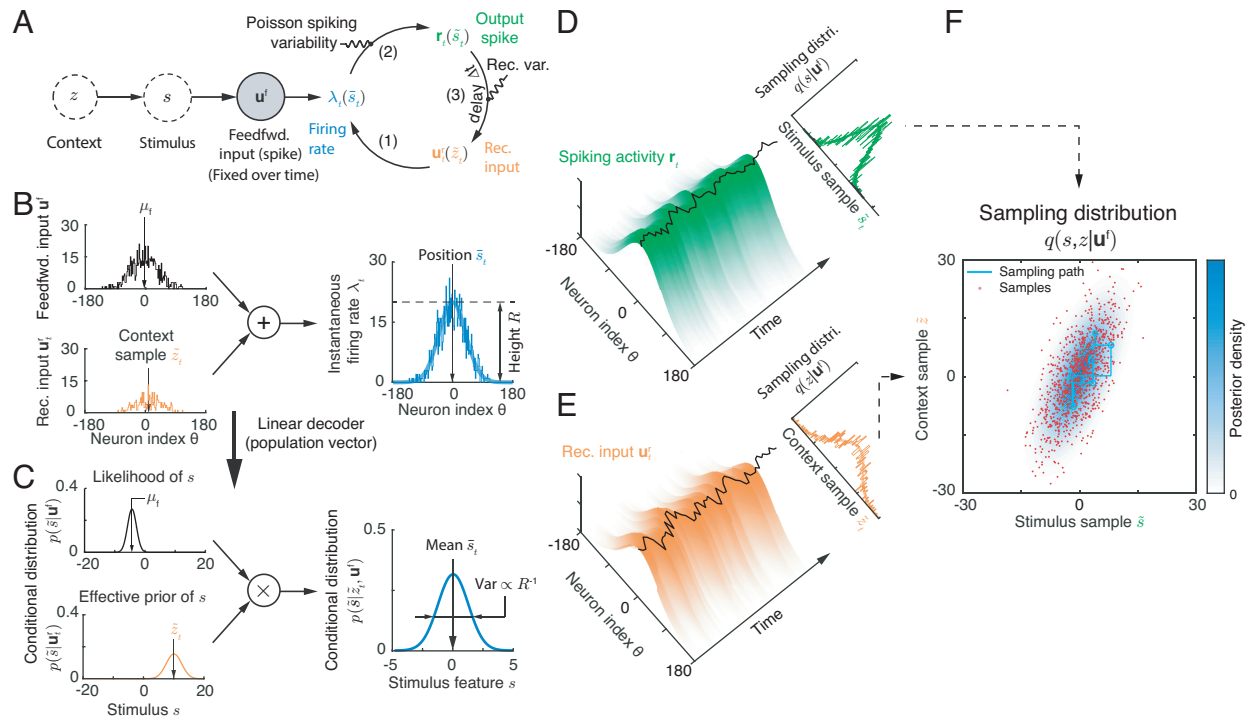


Figure 4: A recurrent circuit generates samples from the posterior defined by a hierarchical generative model. (A) Schematic of recurrent circuit dynamics, in which stimulus, s , and context, z , features are encoded respectively in the population response, \mathbf{r}_t , and recurrent inputs, \mathbf{u}_t^r . (B-C) When the feedforward inputs and recurrent inputs share the same tuning profile, summing the two inputs to define the instantaneous firing rate (B) is equivalent to multiplying the conditional distributions encoded by the two inputs to obtain the conditional distribution of the stimulus, $p(s|z_t, \mathbf{u}^f)$. (C) The conditional distributions of the stimulus can be explicitly read out from corresponding population responses by a linear decoder (B). (D-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 context 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 context (black curve on surface in E) represent samples from the joint posterior of the stimulus and context. The true joint posterior is represented by the blue contour.

220 joint posterior $p(s, z|\mathbf{u}^f)$ (blue contour map in Fig. 3D, Eq. 24) [35]. To approximate $p(s|\mathbf{u}^f)$,
 221 the marginal posterior distribution of s , we can use only samples \tilde{s}_t to obtain the approximating
 222 distribution $q(s|\mathbf{u}^f)$ (compare the two green lines at the margin in Fig. 3D). The same is true for
 223 the marginal posterior over z .

224 Implementing Gibbs sampling of stimulus and context in a recurrently coupled cortical circuit

225 An implementation of Gibbs sampling in a recurrent E circuit can be intuitively understood by
 226 comparing the recurrent network dynamics (Fig. 4A) with the dynamics described by the Gibbs
 227 sampling algorithm (Fig. 3C). In the recurrent network a stimulus sample, \tilde{s}_t , is represented by the
 228 activity of E cells, \mathbf{r}_t , while a context sample, \tilde{z}_t , is represented by recurrent inputs, \mathbf{u}_t^r . To generate

229 correct samples we require that the conditional distribution that is represented by the instantaneous
 230 firing rate, $\boldsymbol{\lambda}_t$ (Eq. 1), matches the conditional distribution used in the Gibbs sampling algorithm
 231 (Eq. 4b), so that $p(\tilde{s}|\tilde{z}_t, \mathbf{u}^f) = p(\tilde{s}|\boldsymbol{\lambda}_t) \propto \exp[\mathbf{h}(\tilde{s})^\top \boldsymbol{\lambda}_t]$. Equating the two distributions (see Eqs. 4a
 232 and 10) yields the relation,

$$\begin{aligned} \ln p(\tilde{s}|\tilde{z}_t, \mathbf{u}^f) &= \ln p(\mathbf{u}^f|\tilde{s}) + \ln p(\tilde{s}|\tilde{z}_t), \\ \Leftrightarrow \mathbf{h}(\tilde{s})^\top \boldsymbol{\lambda}_t &= \mathbf{h}(\tilde{s})^\top \mathbf{u}^f + \mathbf{h}(\tilde{s})^\top \mathbf{u}_t^r. \end{aligned} \quad (5)$$

233 This equation holds when two constraints are satisfied: First, the firing rate vector, $\boldsymbol{\lambda}_t$, needs to
 234 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
 235 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
 236 neural circuit one way for $\boldsymbol{\lambda}_t$ to satisfy these constraints is for feedforward inputs, \mathbf{u}^f , and recurrent
 237 inputs, \mathbf{u}_t^r , to both have Gaussian profiles with the same width, a , as that of $\boldsymbol{\lambda}_t$ (by sharing the same
 238 $\mathbf{h}(\tilde{s})$, Eqs. 5 and 12). This is because the sum of two Gaussian-profile inputs with the same width,
 239 a , gives a firing rate, $\boldsymbol{\lambda}_t$, with the same tuning, as long as the difference of the locations of two
 240 inputs is much smaller than the width, a . Our generative model (Eq. 3) produces feedforward input,
 241 \mathbf{u}^f , with a Gaussian profile and encodes the likelihood function $p(\mathbf{u}^f|\tilde{s})$. The recurrent input, \mathbf{u}_t^r ,
 242 then need to represent the conditional distribution $p(\tilde{s}|\tilde{z}_t)$. Hence, to satisfy Eq. (5) the recurrent
 243 input \mathbf{u}_t^r should have the same Gaussian profile as \mathbf{u}^f (Eq. 29), with its location and magnitude
 244 determined by the mean and precision of $p(\tilde{s}|\tilde{z}_t)$, respectively.

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

251 The recurrent weights must be adjusted so that the recurrent input has the appropriate magni-
 252 tude and width to encode the likelihood $p(s|z)$. To simplify the exposition we first assume that E
 253 neurons are only self-connected, so that the width of recurrent input trivially matches that of the
 254 feedforward input (otherwise recurrence will broaden the profile of the firing rate activity $\boldsymbol{\lambda}_t$ over
 255 the network). To constrain the magnitude of the recurrent weights we require that the sum of the
 256 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
 257 magnitude of the recurrent weights that result in samples from the correct posterior must satisfy:

$$w_E^* = \frac{\langle \mathbf{u}_j^r \rangle}{\langle \mathbf{r}_j \rangle} = \frac{\langle \sum_j \mathbf{u}_j^r \rangle}{\langle \sum_j \mathbf{r}_j \rangle} = \frac{\Lambda_s}{\Lambda_f + \Lambda_s}, \quad (6)$$

258 where Λ_s and Λ_f are the precision of likelihood $p(s|z)$ and $p(\mathbf{u}^f|s)$ respectively (Eq. 3). The optimal

259 recurrent weight, w_E^* , thus encodes the correlation between the stimulus s and the context z . An
260 increase in correlation between s and z , resulting in a narrower diagonal band in $p(s, z)$ (Fig. 3B),
261 requires an increase in the recurrent weight w_E^* for optimal sampling. When context and stimulus
262 are uncorrelated so that $\Lambda_s = 0$, the hierarchical generative model (Fig. 3A) is equivalent to the
263 generative model without context (Fig. 2A) and recurrent interactions are not needed for sampling
264 (i.e., $w_E^* = 0$). Our framework (Eq. 6) thus predicts that optimal Bayesian inference is achieved with
265 recurrent synaptic weights which depend on the correlative structure of the external world. We
266 numerically test this prediction in the next section.

267 **A stochastic E-I spiking network jointly samples stimulus and context**

268 To confirm the predictions of this analysis, we simulated a full recurrent network consisting of both
269 E and I neurons with Poisson spike train statistics (see details in Eqs. 47-50). The E neurons were
270 synaptically connected to each other (Eq. 49, see Fig. 1A), in contrast to the simple network of
271 self-connected E neurons we described above. While recurrent E to E coupling broadens the tuning
272 of excitatory recurrent input, lateral inhibition can sharpen Gaussian firing rate profiles so that it
273 matches that of the feedforward inputs (as required by Eq. 5).

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

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

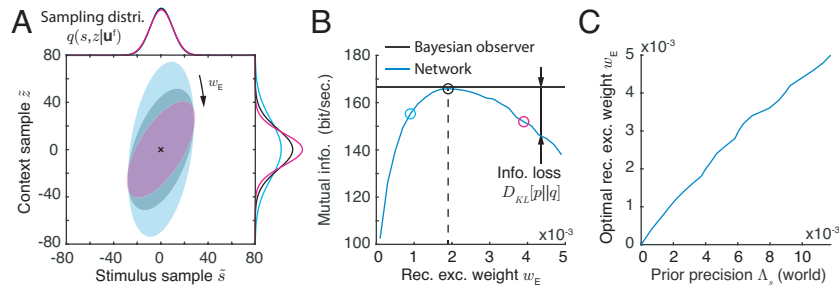


Figure 5: The joint sampling distribution of stimulus and context changes with the recurrent weight in the network. (A) The sampling distribution for different recurrent excitatory weights, w_E . The ratio of excitatory and inhibitory weights was fixed. Ellipses capture three standard deviations from the mean of the joint sampling distribution. Different colors correspond to the three values of w_E , denoted by different symbols in panel B. (B) The mutual information between the latent variables, s and z , 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, $p(s, z | \mathbf{u}^f)$, and the sampling distribution, $q(s, z | \mathbf{u}^f)$. 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 posterior, p (black circle). (C) This optimal weight, w_E^* , increases with prior precision, Λ_s .

294 We use Kullback-Leibler (KL) divergence to measure the distance between the sampling distri-
 295 bution, $q(s, z | \mathbf{u}^f)$, and the true posterior, $p(s, z | \mathbf{u}^f)$ (Eq. 24). The KL divergence quantifies the loss
 296 of mutual information, measured in bits, between the latent variables (s and z) and the feedforward
 297 inputs, \mathbf{u}^f , when the true posterior, p , is approximated by the distribution, q (Eq. 42) [35]. The
 298 mutual information loss in the network is minimized at a unique value of the recurrent weight,
 299 w_E^* , at which the sampling distribution, q , best matches the posterior, p (Fig. 5B, black circle). To
 300 confirm that this optimal recurrent weight, w_E^* , increases with the correlation in the prior (precision
 301 Λ_s , Eq. 6), we numerically obtained the recurrent weight that minimizes the mutual information
 302 loss for each value of Λ_s in the generative model. These results confirmed the predictions of our
 303 theory (Eq. 6, Fig. 5C): When $\Lambda_s = 0$, i.e. when context and stimulus are uncorrelated, a network
 304 with no interactions performs best ($w_E^* = 0$), while for small Λ_s (relative to Λ_f) the optimal weight
 305 w_E^* is positive and increases with Λ_s . In total, we have described a potential mechanism for a
 306 recurrent network of spiking neurons to perform sampling-based Bayesian inference.

307 **Generating samples from multi-dimensional posteriors with coupled neural circuits**

308 To demonstrate the generality of the proposed neural code we next consider a world described
 309 by a broad, rather than deep (hierarchical) generative model. Information about each of two
 310 latent stimuli, $\mathbf{s} = (s_1, s_2)$, is relayed by corresponding feedforward inputs received by a neural
 311 circuit (Fig. 6A). We assume the prior is a bivariate Gaussian distribution (Fig. 6B), i.e., $p(\mathbf{s}) \propto$
 312 $\exp[-\Lambda_s(s_1 - s_2)^2/2] \equiv \mathcal{N}(s_1 - s_2, \Lambda_s^{-1})$, so that Λ_s ($\Lambda_s \geq 0$) characterizes the correlation between s_1
 313 and s_2 . Furthermore, each stimulus, s_m , independently generates feedforward spiking inputs, \mathbf{u}_m^f ,
 314 each of which is received by a separate network and produces responses \mathbf{r}_m for $m = 1, 2$ (Fig. 6A).

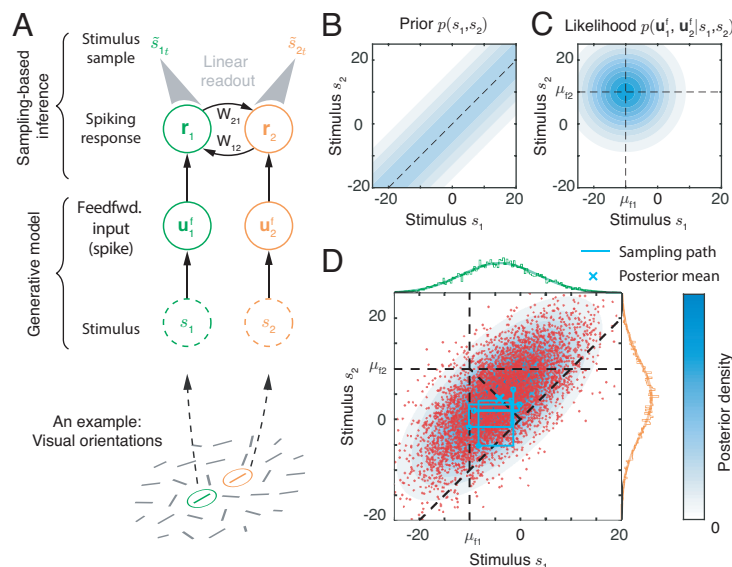


Figure 6: Distributed sampling from a multivariate posterior distributions using coupled networks. (A) Network m ($m = 1, 2$) receives a feedforward input evoked by a stimulus, s_m . The coupling between the two networks represents the stimulus prior. A linear readout from each network, m , can be interpreted as a sample from the posterior of the stimulus, s_m . (B-C) 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 panel (C), $\mu_{f1} = -10$ and $\mu_{f2} = 10$ 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.

315 Thus, the full generative model of the input has the form,

$$\begin{aligned}
 p(\mathbf{u}^f | \mathbf{s}) p(\mathbf{s}) &= \left[\prod_{m=1}^2 p(\mathbf{u}_m^f | s_m) \right] p(s_1, s_2), \\
 &\propto \left[\prod_{m=1}^2 \mathcal{N}(s_m | \mu_{fm}, \Lambda_{fm}^{-1}) \right] \mathcal{N}(s_1 - s_2, \Lambda_s^{-1}).
 \end{aligned}
 \tag{7}$$

316 The likelihood $p(\mathbf{u}_m^f | s_m)$ is the same as that given previously (Eq. 2), where the feedforward inputs,
 317 \mathbf{u}_m^f , are again described by conditionally independent Poisson spike counts with Gaussian tuning
 318 over stimulus s_m . As a concrete example, the two stimuli, s_m , could represent orientations of
 319 local edges falling in the central receptive fields of a V1 hypercolumn (Fig. 6A, bottom), with
 320 each V1 hypercolumn modeled by a network producing the response \mathbf{r}_m (Fig. 6A, top). Then Λ_s
 321 characterizes *a priori* tendency of the stimuli to share similar orientations, and determines how
 322 likely two local edges are to be part of a global line, as in the case of contour integration [42, 43].
 323 However, the generative model defined by Eq. (7) is quite general and has been also used to explain
 324 multisensory cue integration [10] and sensorimotor learning [13].

325 The posterior is a bivariate Gaussian distribution (Fig. 6D, Eq. 34) whose mean is shifted from
 326 the likelihood mean (Fig. 6C) towards to the diagonal line, because of the correlations between

327 the stimuli in the prior (Fig. 6B). We can again use Gibbs sampling to approximate the posterior
 328 $p(\mathbf{s}|\mathbf{u}^f)$ using the following steps,

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

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

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

332 These sampling steps can be implemented distributively in a coupled neural circuit using a
 333 mechanism similar to that we described in the case of a hierarchical generative model. The activity
 334 of each network, \mathbf{r}_m , individually represents samples from the (marginal) posterior of s_m (Fig. 6A,
 335 top). The joint posterior is then approximated as the collection of samples represented by the
 336 activity pairs $(\mathbf{r}_1, \mathbf{r}_2)$. Taking network $m = 1$ as an example, spike response \mathbf{r}_{1t} produces a stim-
 337 ulus sample \tilde{s}_{1t} as long as the instantaneous firing rate λ_{1t} represents the conditional distribution
 338 $p(\tilde{s}_1|\mathbf{u}_1^f, \tilde{s}_{2,t-\Delta t})$ (Eq. 8a). Since the feedforward input, \mathbf{u}_1^f , represents the likelihood $p(\mathbf{u}_1^f|\tilde{s}_1)$, to
 339 obtain the appropriate firing rates, λ_{1t} , the recurrent input from network 2 to network 1, $\mathbf{u}_{12,t}^r$,
 340 must encode the correct conditional distribution, $p(\tilde{s}_{2,t-\Delta t}|\tilde{s}_1)$. As in the case of the mechanism we
 341 proposed to implement sampling as described by Eq. (5), $\mathbf{u}_{12,t}^r$ needs to have the same Gaussian
 342 profile as the firing rate λ_{1t} , the position of $\mathbf{u}_{12,t}^r$ on the stimulus space should match the mean of
 343 $p(\tilde{s}_{2,t-\Delta t}|\tilde{s}_1)$, 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
 344 to the prior correlation, $\Lambda_s \propto \sum_j \mathbf{u}_{12,tj}^r$ (Eq. 39). Hence, each network can sum the feedforward
 345 input and the recurrent input from its counterpart to obtain an update to the instantaneous condi-
 346 tional distribution given by Eq. (8a), and generate independent Poisson spikes to produce a sample
 347 from the instantaneous conditional distribution (Eq. 8b). Notably, the sample of each stimulus
 348 can be locally read out from corresponding network (Eq. 41, Fig. 6A), even if the activities of two
 349 networks are correlated.

350 Since the recurrent input strength represents the stimulus correlation in the prior determined
 351 by precision Λ_s , the coupling between the two networks needs to be tuned to generate the appro-
 352 priate recurrent input. Indeed, in a network with only E neurons, and connections only between
 353 neurons with the same preferred stimulus value but in different networks, the optimal homogeneous
 354 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
 355 obtained with the hierarchical model presented earlier in Eq. (6).

356 **Coupled E-I spiking networks sample bivariate dimensional posteriors**

357 To test the feasibility of the proposed mechanisms for generating samples from a bivariate posterior
 358 we simulated a pair of bidirectionally coupled circuits consisting of E and I neurons (Fig. 7A).

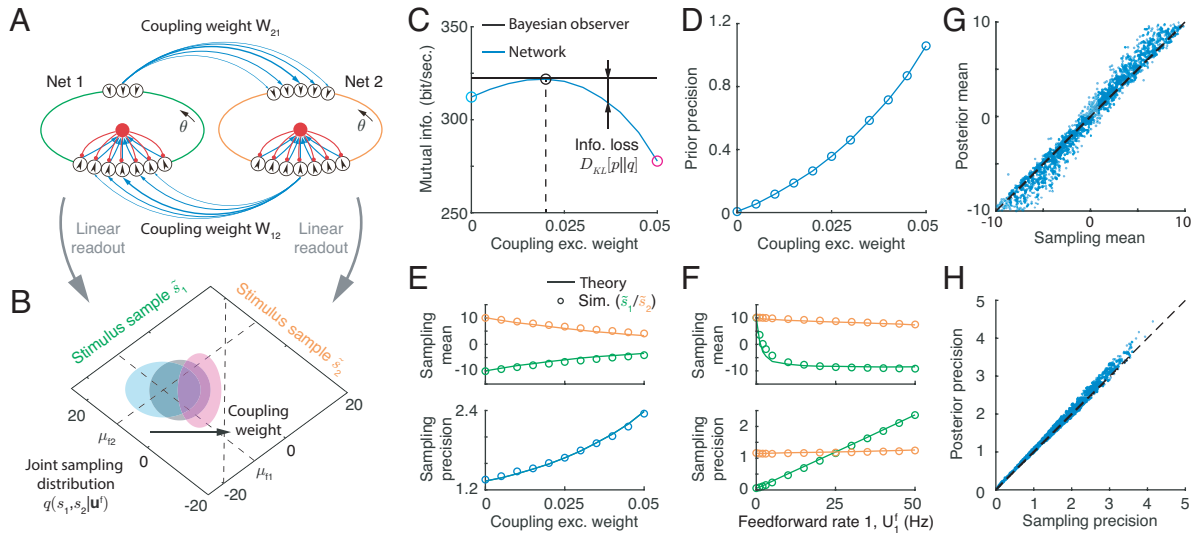


Figure 7: The statistics of the multivariate sampling distribution of stimuli generated by coupled E-I circuits. (A) Each of the two circuits individually generate a sample of a corresponding stimulus which can be read out linearly from that circuit’s activity. Combining the readouts from the two networks 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 panel (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. (G-H) 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.

359 This neural circuit model can be extended to generate samples from higher dimensional posterior
 360 distribution (see Discussion). Each circuit receives feedforward input generated by one of the two
 361 stimuli. On every time step the sample of each stimulus, \tilde{s}_{mt} , can be individually and linearly read
 362 out from the response of corresponding network, \mathbf{r}_{mt} (Eq. 41). Jointly, the two stimulus samples,
 363 one each from both networks, $\tilde{\mathbf{s}}_t = (\tilde{s}_{1t}, \tilde{s}_{2t})^\top$, provide a sample from the joint posterior of the
 364 two latent stimuli (Fig. 7B). We assumed that the synaptic connections between the networks,
 365 w_{mn} ($m, n = 1, 2; m \neq n$), are excitatory, but target both E and I neurons, while inhibitory
 366 connections are local to each network. We also adjusted network parameters so that the profiles
 367 of the inputs across networks (e.g., the inputs from network 2 to 1) have the same tuning profile
 368 as the feedforward inputs (see Methods). Since we assumed uniform marginal priors (see Eq. 32),

369 recurrent connections between E neurons within the a circuit were absent, while E and I neurons
370 within a circuit were recurrently connected to ensure network stability. For simplicity, we chose
371 parameters so that the two circuits were symmetric, but the strength of the feedforward inputs to
372 each could differ.

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

385 We next asked how the sampling distribution in the network depends on network and feedfor-
386 ward input parameters. As the coupling between the two circuits increased, the sample means of
387 both stimuli converge (Fig. 7E, top) and the sampling precision of both stimuli increased as well
388 (Fig. 7E, bottom), in agreement with a more correlated stimulus prior. We also tested whether
389 a network with fixed parameters can generate samples from a family of posteriors with different
390 uncertainties. To do so, we changed the uncertainty of the likelihood of s_1 by changing the fir-
391 ing rate in the feedforward input \mathbf{u}_1^f received by network 1. We observed that with a narrower
392 likelihood of s_1 , the sample means of both stimuli shifted towards the mean of likelihood of s_1
393 (-10°), and sampling precision increased, consistent with a change in the posterior distribution
394 (Fig. 7F). Lastly, to demonstrate the robustness of this network implementation of sampling-based
395 inference we compare the sampling distributions to the true posteriors under different combinations
396 of input and network parameters (Fig. 7G-H), in each case setting the recurrent coupling to the
397 optimal value, w_{mn}^* , obtained numerically. Across different parameter values we observe excellent
398 agreement in both the mean (Fig. 7G) and precision (Fig. 7H) of the two densities. In sum, our
399 recurrent network of spiking neuron models can be extended to support sampling-based Bayesian
400 inference with multi-dimensional stimuli.

401 **A signature of stimulus sampling: internally generated differential noise correlations**

402 A central prediction of our circuit framework for sampling-based Bayesian inference is that an
403 increase in the correlation between stimuli in the sensory world should result in stronger synapses

404 between neurons whose activities represent these stimuli (see Eq. 6). This is a difficult prediction
 405 to test since measuring synaptic connectivity along a functional axis is already challenging [44], let
 406 alone measuring a change in synaptic strength owing to a change in stimulus statistics. Here we
 407 outline a testable prediction of our theory by identifying a measurable, population-level signature
 408 of changes in functionally related recurrent synaptic strengths.

409 In response to a fixed feedforward input the responses of a recurrent circuit implementing stim-
 410 ulus sampling will fluctuate. The alignment of the recurrent circuitry and neuronal stimulus tuning
 411 causes a portion of these activity fluctuations to align with the subspace in which stimuli are coded.
 412 As an example, consider the sampling implemented by a single recurrent network (Fig. 4A), and
 413 suppose the population response fluctuates around its mean position (0° in the example of Fig. 8A),
 414 ignoring fluctuations along other directions in neuronal response space. The activity of neuron pairs
 415 with stimulus preference both above or below the mean position are positively correlated (the black
 416 and blue neurons in Fig. 8A), while the activity of neuron pairs with preferences straddling the mean
 417 are negatively correlated (the black and red neurons in Fig. 8A). Such stimulus sampling generates
 418 a covariance component which is proportional to the outer product of the derivative of neuronal
 419 tuning (Fig. 8B), i.e., $\mathbf{f}'_s \mathbf{f}'_s{}^\top$, where \mathbf{f}'_s denotes the derivative of tuning $\mathbf{f}(s) = \langle \boldsymbol{\lambda}_t \rangle$ (mean firing rate)
 420 over stimulus s . Such noise correlations have been referred to as *differential correlations* [4, 17], and
 421 are generally viewed as deleterious to stimulus coding. Stochastic sampling in coupled networks
 422 (Fig. 6A) produces similar differential noise correlations (see Supplemental Information).

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

$$\Sigma_{DC} = V(\bar{s}|\mathbf{u}^f) \mathbf{f}'_s \mathbf{f}'_s{}^\top, \quad (9)$$

$$\text{where } V(\bar{s}|\mathbf{u}^f) = \frac{\Lambda_s}{\Lambda_f(\Lambda_f + \Lambda_s)} = a^2 n_f^{-1} w_E^*,$$

426 where $V(\bar{s}|\mathbf{u}^f)$ is the variance of \bar{s}_t in equilibrium over time, and \bar{s}_t is the mean of the instan-
 427 taneous conditional distribution (Eq. 4a) represented by the position of instantaneous firing rate
 428 $\boldsymbol{\lambda}_t$ (Fig. 2B). Importantly, the amplitude of differential correlations increases with the recurrent
 429 weight, w_E^* , which is set by the prior precision Λ_s (Eq. 6; Fig. 8C). Thus, in our framework inter-
 430 nally generated differential correlations are a by-product of inference by sampling from posterior
 431 distributions of stimuli in a structured world.

432 Distinguishing external and internal differential correlations

433 The previous analysis of internally generated differential correlations in a circuit implementing
 434 sampling-based inference is based on the assumption of a fixed feedforward input (Eq. 9). However,

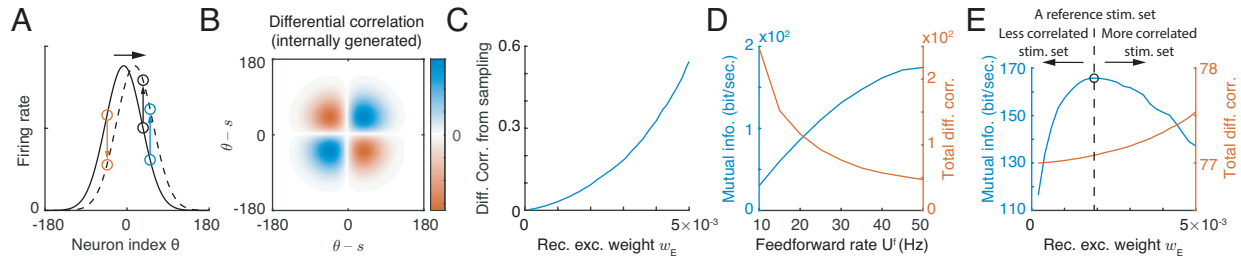


Figure 8: Stimulus sampling by a network is reflected in the internally generated differential correlations, whose impact differs from differential correlations inherited from feedforward inputs. (A) Stimulus sampling via spike generation causes the population firing rate to fluctuate along the stimulus 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.

435 in typical neurophysiology experiments an external stimulus, s , is fixed, while the feedforward
 436 input, \mathbf{u}^f , fluctuates due to variability in sensory acquisition and transmission noise (Eqs. 3 and 7).
 437 Hence, differential correlations of neuronal population responses are a combination of correlations
 438 inherited from feedforward input [45], and correlations generated by recurrent network interactions
 439 that align with the population stimulus tuning [24]. When the feedforward input is described by a
 440 hierarchical generative model (Eq. 2), the total magnitude of differential correlations in the evoked
 441 response is $a^2 n_f^{-1} w_E \mathbf{f}'_s \mathbf{f}'_s{}^\top + a^2 n_f^{-1} \mathbf{f}'_s \mathbf{f}'_s{}^\top$ (see Eq. 46), where the second term reflects differential
 442 correlations inherited from the feedforward input (compare with Eq. 9). Although the two sources
 443 of differential correlations are intertwined in the neuronal response, they impact the information
 444 content differently thus offering a potential way to distinguish between them in neural data.

445 Externally generated differential correlations decrease with feedforward input rate which could
 446 be modulated by visual stimulus strength such as contrast (Fig. 8D, red curve). As a consequence,
 447 the mutual information (the information between feedforward inputs \mathbf{u}^f and the latent variables, i.e.,
 448 s and z , sampled by recurrent network in Fig. 4A, Eq. 42) increases with feedforward input intensity
 449 (Fig. 8A, blue curve). We therefore have a monotonic, decreasing relationship between externally
 450 generated differential correlations and mutual information. This is expected since such inherited
 451 correlations always impair information processing, as observed previously [4, 17]. In contrast, an
 452 increase in recurrent weights, w_E , increases internally generated differential correlations, but results
 453 in a non-monotonic change in mutual information (Fig. 8B). Hence there is a non-monotonic relation
 454 between internally generated differential correlations and the mutual information between stimulus
 455 and feedforward inputs. In sum, the impact of external and internal differential correlations on
 456 stimulus coding can be distinguished by their respective monotonic and non-monotonic relation

457 with the mutual information between stimulus and response.

458 Discussion

459 We have presented a framework in which neuronal response variability and recurrent synaptic con-
460 nections, two ubiquitous features of cortex, are jointly used to implement sampling-based Bayesian
461 inference in neuronal circuit models. Combining mathematical analysis and network simulations we
462 established that stereotypical Poisson variability of discrete spike counts can drive flexible sampling
463 from a family of continuous distributions. The sampling statistics are determined by the structure
464 of recurrent coupling, which stores information about the stimulus prior, and feedforward inputs
465 which convey the stimulus likelihood. Sampling-based inference is implemented in two steps: the
466 instantaneous firing rate, determined by the sum of feedforward and recurrent inputs, represents
467 the instantaneous conditional distribution of latent stimulus, while Poissonian variability in spike
468 generation is used to generate a random stimulus sample from this conditional distribution. A sim-
469 ple circuit model is able to generate samples from multi-dimensional posteriors of latent variables
470 organized hierarchically or in parallel, which underlies the computational basis of a wide range of
471 perceptual and cognitive processes [46].

472 Comparison with other neural coding frameworks

473 The neural code we described shares some features with codes described in previous studies, includ-
474 ing parametric representations in probabilistic population codes (PPCs) [15, 36, 37], and sampling-
475 based codes (SBCs) [16, 27–32]. In our framework the conditional distributions of latent variables
476 is represented by instantaneous firing rates which linearly encode the logarithms of these conditional
477 distributions, and have a mathematical form that is similar to that used in past studies describing
478 PPCs (e.g., Eq. 5). Further, the posterior is represented by stimulus samples generated through a
479 random process, a feature of all SBCs. Despite these similarities, there are fundamental differences
480 between the neural code we described and previously proposed PPCs and SBCs.

481 PPCs are generally implemented in networks with no internally generated variability, with
482 stochasticity inherited from the stimulus. In contrast, our proposed network is doubly stochastic:
483 The Poisson variability in the feedforward input allows a single realization of the feedforward input
484 to represent the whole stimulus likelihood [36], while internally generated Poisson variability drives
485 stimulus sampling. Further, in PPCs the posterior is represented parametrically by a one-shot
486 neuronal response, while in our proposed network the joint posterior is approximated by a sequence
487 of samples, each obtained as a linear readout from the instantaneous neuronal responses. Although
488 it takes time to collect sufficient samples to approximate the posterior, a computational benefit
489 compared with PPCs is that inference of a multivariate posterior can be implemented by linearly
490 coupled networks (Fig. 6), while in PPCs nonlinear coupling between networks is required [47].

491 Conventional SBCs are used to generate samples directly in a neural space whose dimension is
492 given by the number of neurons in the population [16, 27, 28, 30–33], where a neuronal response,
493 \mathbf{r}_t , is interpreted directly as a sample from the (marginal) posterior of neuronal responses, $p(\mathbf{r})$.
494 Hence the posterior mean is the temporally averaged population response, and the covariance of
495 population responses is the posterior covariance. In contrast, our proposed network generates sam-
496 ples in a low dimensional stimulus subspace embedded in high dimensional neural activity space.
497 The linear projection of network activity, \mathbf{r}_t , onto the stimulus subspace represents a sample from
498 the stimulus posterior, similar to a previous study [29]. A computational benefit of sampling in a
499 low dimensional stimulus subspace is convergence speed, as the volume of the stimulus subspace is
500 significantly smaller than that of the neural activity space. Indeed, in our examples sequences of
501 samples generated by a single recurrent network (Fig. 4) and coupled networks (Fig. 6) can both
502 converge to an equilibrium distribution in less than 20ms, which is fast enough to complete inference
503 on a behaviorally relevant time scale (Fig. S6). Furthermore, the multiplication of probability dis-
504 tributions of latent stimulus, which is central to Bayesian inference (e.g., cue combination, decision
505 making, see review in [15]), can be implemented by summing the inputs to a neuronal population
506 (Eq. 5). This follows from the fact that the instantaneous population input (or firing rate) linearly
507 encodes the logarithm of a probability distribution (Eqs. 1 and 5). In contrast, producing samples
508 in neural activity space using conventional SBCs requires nonlinear operations in neural circuits in
509 order to multiply probability distributions (or histograms) of the samples [15].

510 A recent study demonstrated that an E-I recurrent network of rate-based neurons can be nu-
511 merically optimized for sampling-based Bayesian inference [32]. In contrast, we used a theoretical
512 approach to derive a network model of simplified spiking neurons which implements sampling-based
513 inference. This allowed us to explicitly describe the putative neural mechanisms needed for such
514 sampling. Although the two studies use different generative models and neural representations,
515 the network models in both studies share some common characteristics: ring structure, Poisson-
516 like response variability, and tuning-dependent noise correlation (Fig. S1D). This implies that the
517 seemingly different generative models and neural representations in the two studies reflect more
518 general principles, as suggested in [48]. It will be interesting to extend our theoretical approach
519 to dynamical spiking neurons to determine how the timescales of neuronal dynamics and neuronal
520 oscillations impact inference in rich, dynamic sensory scenes (see below).

521 **Testing the prediction that recurrent synaptic strength is determined by correlations** 522 **between latent stimuli**

523 Differential noise correlations generated by recurrent network interactions are a signature of network
524 sampling in our framework (Fig. 5C and 8C). This is in contrast to earlier studies where differential
525 correlations were inherited from feedforward inputs [17, 49]. While internally generated differential

526 correlations could also emerge from a recurrent circuit which is not implementing inference [22, 24,
527 49–52] or implementing inference via other algorithms [53], in our framework the relation between
528 the magnitude of internally generated differential correlations, the posterior uncertainty, and the
529 strength of the recurrent synaptic weight (Eq. 9) provides a clear test which can be used to verify
530 our proposed circuit mechanism of sampling-based inference. One possible experimental approach
531 would modulate the functional recurrent strength by using a perceptual learning task. Specifically,
532 after using a reference stimulus set with a prescribed correlation between latent stimuli to fully
533 train an animal, we expect that recurrent synaptic weights will strengthen or weaken to improve
534 inference (Fig. 8E, dashed line). This will result in a fixed value of differential noise correlations
535 in the population response due to the recurrent circuitry. Re-training with a stimulus set that has
536 more (less) correlated latent stimuli compared to the reference set will cause the recurrent weights
537 to increase (decrease) (Fig. 8E, red line). When the reference stimulus set is again used to drive task
538 behavior, then performance (as a proxy of mutual information) will decrease, regardless of whether
539 differential correlations have increased or decreased compared to those resulting from the reference
540 stimulus set (Fig. 8E, arrows). In brief, the non-monotonic relationship between differential noise
541 correlations and the mutual information between stimulus and responses which support Bayesian
542 inference offers a clear (and falsifiable) experimental prediction.

543 **Extensions of circuit-based Bayesian inference**

544 Implementing sampling-based inference in our proposed network requires that feedforward and re-
545 current inputs have the same tuning profile over the stimulus (Eq. 5). This assumption is supported
546 by experiments in layers 4 and 2/3 in mouse V1 [8]. Moreover, the recurrent connections in our
547 network model are translation-invariant in the stimulus subspace, an assumption widely used in
548 continuous attractor networks (CAN) [22, 51, 54, 55]. Translation-invariant connections simplify
549 the mathematical analysis, but are not required for a circuit to implement sampling. Adding ran-
550 domness in recurrent connectivity only increases the variance of the sampling distribution. In the
551 past, CANs have been shown to achieve maximal likelihood estimation (point estimate) via template
552 matching [15, 55, 56]. Here we have shown that a network with CAN-like structure and internally
553 Poisson spiking variability is able to perform sampling-based Bayesian inference. In our network
554 correlations in the stimulus prior are represented by the strength of recurrent synaptic activity,
555 which implies that the (subjective) prior precision in the network increases with the feedforward
556 input strength. To maintain a fixed prior in the network recurrent weights need to decrease with
557 increased feedforward input strength which encodes the likelihood precision, Λ_f (Eq. 6). There-
558 fore, the (subjective) prior stored in the network with fixed recurrent weights may differ from the
559 objective stimulus prior in the world (Λ_s in Eqs. 3 and 7) with feedforward inputs of different
560 strengths. This could be solved by short-term synaptic depression which decreases the synaptic

561 efficacy at increased neuronal firing rates [57]. On the other hand, since the proposed recurrent
562 circuit is general, this result may explain the origin of inductive bias [58] or confirmation bias [59]
563 in cortical processing. Another possibility is that the recurrent circuit represents a more complex
564 generative model which better captures the statistical structure of natural stimuli [30, 32, 60]. We
565 only considered sampling driven by spiking variability with a Fano factor of 1, while cortical re-
566 sponses often have Fano factors that differ from 1 [61, 62]. In the latter case, our theory can still
567 work by changing the feedforward connection weight to compensate for the change in Fano factor,
568 as suggested in a recent study [63].

569 To keep our exposition transparent we only presented models with minimal complexity. Our
570 proposed network mechanism of sampling-based inference can be generalized to more complex gen-
571 erative models, since the assumption of Gaussianity (Eqs. 21 and 22) and the analytical expression
572 in Eq. (24) are not essential, and several relaxed frameworks may be explored. First, similar
573 networks can generate samples from other multi-dimensional distributions where the conditional
574 distribution of each latent variable belongs to the linear exponential family [35, 36]. This could be
575 done by changing the tuning functions of neurons to another appropriate profile, as the logarithm
576 of tuning determines the type of sampling distribution (Eq. 1). When sampling from non-Gaussian
577 distributions, the stimulus samples can be linearly read out with the weight determined by the
578 tuning profile (i.e., $\mathbf{h}(s)$ in Eq. 1, [36]). Second, the tuning of recurrent inputs does not need to
579 be the same as that of feedforward inputs. Instead the logarithm of recurrent input tuning can
580 have a form of the conjugate prior with the likelihood conveyed by feedforward inputs. Third, the
581 network model could also be used to infer the latent variables with a non-uniform marginal prior,
582 if, for example, the preferred stimuli of neurons in the population are not distributed uniformly
583 in the stimulus subspace [64]. Lastly, we considered only non-structured inhibition for simplic-
584 ity. Structured inhibitory connections could modulate the position of excitatory responses in the
585 stimulus subspace, i.e., the mean of the conditional distribution. Such interplay between E and I
586 neurons with structured inhibition has the potential to implement Hamiltonian sampling, where
587 the I neurons represent the sample of auxiliary variables [33, 35].

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

593 Methods

594 A linear network of excitatory neurons

595 We study how a generic recurrent network model consisting solely of N_E excitatory (E) neurons
 596 with Poisson spiking statistics (no inhibitory neurons) can implement sampling-based Bayesian
 597 inference to approximate the stimulus posterior. We describe neuronal activity using a time-
 598 discretized Hawkes process (a type of multivariate, inhomogeneous Poisson process [65]). The
 599 instantaneous firing rates of the neurons in the network at time t , λ_t , obey the following recurrent
 600 equations:

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

$$\mathbf{r}_t \sim \prod_{j=1}^{N_E} \text{Poisson}(\lambda_{tj} \Delta t), \quad (11)$$

601 where \mathbf{u}^f is the feedforward Poisson spiking input (described below; Eq. 18), \mathbf{u}_t^r is the continuous
 602 valued recurrent input at time t , and $\boldsymbol{\xi}_t$ is a N_E dimensional independent Gaussian white noise.
 603 Hence, over each time interval $[t - \Delta t, t]$ the activity of the neurons in the network is modeled by
 604 a vector of independently generated Poisson spike counts, \mathbf{r}_t , with means determined by the rates
 605 λ_t . The parameters w_E and σ_r determine the excitatory recurrent weight and recurrent variability,
 606 respectively.

607 Poisson spike generation samples stimulus

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

$$\lambda_{tj} = R \exp[-(\bar{s}_t - \theta_j)^2 / 2a^2] = R \exp[\mathbf{h}_j(\bar{s}_t)], \quad (12)$$

613 where \bar{s}_t characterizes the position of the population firing rate on the stimulus subspace (Fig. 1B,
 614 x-axis), while R and a denote the height and width of the population firing rate, respectively.
 615 Further, θ_j is the preferred stimulus value of neuron j , and the preferred stimuli of all neurons,
 616 $\{\theta_j\}_{j=1}^{N_E}$, are uniformly distributed over the range of stimulus s (Fig. 1B).

617 To simplify the analysis, we first assume that the instantaneous firing rate is fixed over time.
 618 When generating Poisson spikes \mathbf{r}_t from λ_t , the probability of observing a stimulus sample \bar{s}_t

619 (embedded in \mathbf{r}_t) can be derived as (see details in Supplemental Information),

$$\begin{aligned} p(\mathbf{r}_t|\boldsymbol{\lambda}_t) &= \prod_{j=1}^{N_E} \text{Poisson}(\mathbf{r}_{tj}|\boldsymbol{\lambda}_{tj}\Delta t), \\ &\propto \exp[\mathbf{h}(\bar{s}_t)^\top \mathbf{r}] \cdot [n_{\boldsymbol{\lambda}}^{n_{\mathbf{r}}} \exp(-n_{\boldsymbol{\lambda}})], \\ &\propto \mathcal{N}(\tilde{s}_t|\bar{s}_t, a^2 n_{\mathbf{r}}^{-1}) \text{Poisson}(n_{\mathbf{r}}|n_{\boldsymbol{\lambda}}), \end{aligned} \quad (13)$$

620 where $n_{\mathbf{r}} = \sum_j \mathbf{r}_{tj}$ is the number of emitted spikes across the whole neural population, and $n_{\boldsymbol{\lambda}} =$
621 $\sum_j \langle \boldsymbol{\lambda}_j \rangle \Delta t$ is the sum of population firing rate. Here $\mathcal{N}(s|\mu, \sigma^2)$ denotes a Gaussian distribution
622 with mean μ and variance σ^2 , and $\mathbf{h}(\bar{s}_t)$ is a vector with the j^{th} element as $\mathbf{h}_j(\bar{s}_t)$ shown in Eq. (12).
623 The logarithm of the firing rate profile, $\mathbf{h}(\bar{s}_t)$, determines how the stimulus sample \tilde{s}_t and its mean,
624 \bar{s}_t , can be read out respectively from \mathbf{r}_t and $\boldsymbol{\lambda}_t$,

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

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

626 The sampling variability of \tilde{s}_t in a single time step depends on the number of emitted spikes,
627 $n_{\mathbf{r}}$. When the fixed rates, $\boldsymbol{\lambda}_t$, repeatedly generate spikes over time, the sampling distribution of
628 \tilde{s}_t can be calculated by marginalizing the likelihood (Eq. 13, last line) over different values of $n_{\mathbf{r}}$
629 since $n_{\mathbf{r}}$ varies across time (detailed calculation by using Laplacian approximation can be seen in
630 Supplemental Information),

$$\begin{aligned} p(\tilde{s}_t|\boldsymbol{\lambda}_t) &= \sum_{n_{\mathbf{r}}} \mathcal{N}(\tilde{s}_t|\bar{s}_t, a^2 n_{\mathbf{r}}^{-1}) \text{Poisson}(n_{\mathbf{r}}|n_{\boldsymbol{\lambda}}), \\ &\approx \mathcal{N}(\tilde{s}_t|\bar{s}_t, a^2 n_{\boldsymbol{\lambda}}^{-1}). \end{aligned} \quad (15)$$

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

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

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

636 Eq. (16) suggests that the type of sampling distribution (or the conditional distribution) that
637 is obtained from spike generation variability is determined by the profile of the instantaneous firing
638 rate, i.e., $\mathbf{h}(\bar{s}_t)$ (Eq. 12). Although the sampling distribution belongs to the linear exponential
639 family of distributions which is similar with the probabilistic population code (PPC) [36], there are
640 different ways in representing these distributions. In PPCs the likelihood over \bar{s}_t is parametrically
641 represented by a single realization of independent neuronal response \mathbf{r} (Eq. 13), while in our work

642 the distribution is approximated by a sequence of samples, \tilde{s}_t , effectively generated by conditionally
 643 independent Poisson spike discharges.

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

$$\lambda_{tj} = R_t \exp[\mathbf{h}_j(\bar{s}_t)] + \delta_{\perp} \lambda_{tj}, \quad (17)$$

647 where $\delta_{\perp} \lambda_t$ denotes the deviation from a smooth Gaussian profile. Note that the sampling dis-
 648 tribution only depends on the position, \bar{s}_t , and the sum of instantaneous firing rate, n_{λ} (Eq. 16),
 649 which corresponds to two perpendicular directions in the N_E dimensional space of λ_t . For any
 650 instantaneous firing rate vector, λ_t , we can always find \bar{s}_t and R_t that make the deviation $\delta_{\perp} \lambda_t$
 651 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
 652 implies that deviations from Gaussian firing rate profiles do not affect our theory.

653 Feedforward spiking input conveys the likelihood of stimulus

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

$$\begin{aligned} p(\mathbf{u}^f | s) &= \prod_{j=1}^{N_E} \text{Poisson}[\mathbf{u}_j^f | \langle \mathbf{u}_j^f(s) \rangle], \\ \langle \mathbf{u}_j^f(s) \rangle &= U^f \exp[\mathbf{h}_j(s)] = U^f \exp[-(\theta_j - s)^2 / 2a^2]. \end{aligned} \quad (18)$$

656 Here \mathbf{u}_j^f denotes the feedforward input received by the j^{th} E neuron, and $\langle \mathbf{u}_j^f(s) \rangle$ is the tuning of
 657 the feedforward input. This mathematical description of feedforward input is the same as the one
 658 used in the definition of typical PPCs [15, 36, 37]. Since the preferred stimulus values, $\{\theta_j\}_{j=1}^{N_E}$, of
 659 all feedforward inputs are uniformly distributed in stimulus space then the likelihood of s given a
 660 single observation of the input, \mathbf{u}^f , satisfies [36, 37],

$$\begin{aligned} p(\mathbf{u}^f | s) &\propto \exp[\mathbf{h}(s)^{\top} \mathbf{u}^f], \\ &\propto \mathcal{N}(s | \mu_f, \Lambda_f^{-1}). \end{aligned} \quad (19)$$

661 The logarithm of tuning, $\mathbf{h}(s)$, determines the type of likelihood [15]. Specifically, the Gaussian
 662 tuning leads to a Gaussian likelihood (Eq. 19), whose mean, μ_f , and precision, Λ_f , are both linear
 663 functions of the inputs,

$$\mu_f = n_f^{-1} \sum_j \mathbf{u}_j^f \theta_j, \quad \Lambda_f = a^{-2} n_f = a^{-2} \sum_j \mathbf{u}_j^f. \quad (20)$$

664 The mean, μ_f , represents the position of \mathbf{u}^f in stimulus subspace, and the precision, Λ_f , is propor-
 665 tional to the sum of total feedforward spike counts, n_f .

666 **A recurrent network samples hierarchical latent variables**

667 **A hierarchical generative model**

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

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

672 where $\mathcal{U}(a, b)$ denotes a uniform distribution over $[a, b]$. The assumption of a uniform prior, $p(z)$,
 673 simplifies our model significantly, as it implies the spatial homogeneity of the network model as
 674 given by Eqs. (18-19). However, this assumption is not essential for our main results. Due to
 675 the differences between the stimulus (local) and context (global) aspects of the sensory scene, the
 676 stimulus, s , is not identical to the context z , but we assume that the two are correlated, so that

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

677 In sum, the whole generative model is determined by,

$$\begin{aligned} p(\mathbf{u}^f, s, z) &= p(\mathbf{u}^f|s)p(s|z)p(z), \\ &\propto \mathcal{N}(s|\mu_f, \Lambda_f^{-1})\mathcal{N}(s|z, \Lambda_s^{-1}), \end{aligned} \quad (23)$$

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

679 **Approximate Bayesian inference via Gibbs sampling**

680 The joint posterior of stimulus and context can be analytically derived given the generative model
 681 (Eq. 23),

$$\begin{aligned} p(s, z|\mathbf{u}^f) &= \mathcal{N}[(s, z)^\top | \boldsymbol{\mu}_p, \mathbf{K}_p^{-1}], \\ \boldsymbol{\mu}_p &= (\mu_f, \mu_f)^\top, \quad \mathbf{K}_p = \begin{pmatrix} \Lambda_f + \Lambda_s & -\Lambda_s \\ -\Lambda_s & \Lambda_s \end{pmatrix}. \end{aligned} \quad (24)$$

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

684 We use the stochastic response of our recurrent network (Eqs. 10-11), as a basis for Gibbs
 685 sampling [31, 35, 41] (a type of Monte Carlo method) to approximate the joint posterior of stimulus,
 686 s , and context, z . To describe the iterative Gibbs algorithm, we assume that a context sample, \tilde{z}_t ,
 687 is provided at time t , which is then combined with the feedforward input to update the conditional
 688 distribution of stimulus s (step 1 in Fig. 3C),

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

$$\bar{s}_t = \frac{\Lambda_f \mu_f + \Lambda_s \tilde{z}_t}{\Lambda_f + \Lambda_s}, \quad \Lambda = \Lambda_f + \Lambda_s. \quad (25)$$

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

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

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

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

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

696 Implementing the Gibbs sampling in a recurrent circuit model

697 Gibbs sampling of the stimulus (Eq. 4b) can be implemented via independent Poisson spike gener-
 698 ation, as long as the conditional distribution encoded in $\boldsymbol{\lambda}_t$ (Eq. 16) is the same as the conditional
 699 distribution in the Gibbs sampling algorithm (Eq. 4a), i.e., $\ln p(\tilde{s}|\boldsymbol{\lambda}_t) = \mathbf{h}(\tilde{s})^\top \boldsymbol{\lambda}_t = \ln p(\tilde{s}|\tilde{z}_t, \mathbf{u}^f)$.
 700 This condition can be realized in the recurrent circuit by relating the expressions describing the
 701 neural dynamics (Eq. 10) and those describing the Gibbs sampling distribution (Eq. 4a) to yield,

$$\begin{aligned} \ln p(\tilde{s}|\tilde{z}_t, \mathbf{u}^f) &= \mathbf{h}(\tilde{s})^\top \boldsymbol{\lambda}_t, \\ &= \mathbf{h}(\tilde{s})^\top \mathbf{u}^f + \mathbf{h}(\tilde{s})^\top \mathbf{u}_t^r, \\ &= \ln p(\mathbf{u}^f|\tilde{s}) + \ln p(\tilde{s}|\tilde{z}_t). \end{aligned} \quad (27)$$

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

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

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

$$\begin{aligned}\mathbf{u}_{tj}^r(\tilde{z}_t) &= U^r \exp[-(\theta_j - \tilde{z}_t)^2/2a^2] + \delta_{\perp} \mathbf{u}_{tj}^r, \\ \tilde{z}_t &= \sum_j \mathbf{u}_{tj}^r \theta_j / \sum_j \mathbf{u}_{tj}^r, \quad \Lambda_s = a^{-2} \sum_j \mathbf{u}_{tj}^r,\end{aligned}\tag{29}$$

705 whose position on the stimulus subspace is \tilde{z}_t , and the sum of input (height) is determined by Λ_s ,
706 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
707 deviation from a smooth Gaussian and is perpendicular to the direction of \tilde{z}_t and Λ_s .

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

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

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

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

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

719 where $[\cdot]_+$ denotes negative rectification. Here $\boldsymbol{\xi}_t$ is a N_E dimensional Gaussian white noise with
720 $\langle \boldsymbol{\xi}_t(i) \boldsymbol{\xi}_{t'}(j) \rangle = \delta_{ij} \delta(t - t')$, δ_{ij} and $\delta(t - t')$ are Kronecker and Dirac delta functions respectively, $\bar{\mathbf{u}}_t^r$
721 represents the conditional distribution $p(\tilde{z}|\tilde{s}_{t-\Delta t})$, and \mathbf{u}_t^r represent a context sample \tilde{z}_t (Eq. 29).
722 The multiplicative variability on recurrent interaction may come from synaptic noise [66, 67].

723 **Coupled circuits sample a multi-dimensional posterior**

724 We consider a generative model which has multiple latent stimuli, $\mathbf{s} = (s_1, s_2, \dots, s_m)$, which are
725 organized in parallel (Fig. 6A). Without loss of generality, we consider the simplest case where
726 $m = 2$, and the same mechanism can be straightforwardly extended to any $m > 2$. We assume the

727 joint prior of \mathbf{s} is a multivariate normal distribution,

$$p(\mathbf{s}) = \mathcal{N}(\mathbf{s}|\boldsymbol{\mu}_s, \boldsymbol{\Lambda}_s^{-1}) \propto \exp[-\Lambda_s(s_1 - s_2)^2/2],$$

$$\text{with } \boldsymbol{\Lambda}_s = \Lambda_s \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}, \quad (32)$$

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

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

$$p(\mathbf{u}^f|\mathbf{s})p(\mathbf{s}) = \left[\prod_{m=1}^2 p(\mathbf{u}_m^f|s_m) \right] p(s_1, s_2),$$

$$\propto \left[\prod_{m=1}^2 \mathcal{N}(s_m|\mu_{fm}, \Lambda_{fm}^{-1}) \right] \mathcal{N}(\mathbf{s}|\boldsymbol{\mu}_s, \boldsymbol{\Lambda}_s^{-1}), \quad (33)$$

$$\propto \mathcal{N}(\mathbf{s}|\boldsymbol{\mu}_f, \boldsymbol{\Lambda}_f^{-1}) \mathcal{N}(\mathbf{s}|\boldsymbol{\mu}_s, \boldsymbol{\Lambda}_s^{-1}),$$

741 where $\boldsymbol{\mu}_f = (\mu_{f1}, \mu_{f2})^\top$, and the likelihood precision matrix $\boldsymbol{\Lambda}_f = \text{diag}(\Lambda_{f1}, \Lambda_{f2})$ is a diagonal matrix.

742 Gibbs sampling of the multi-dimensional posterior in a coupled neural circuit

743 Given the generative model (Eq. 33), the joint posterior of s_1 and s_2 is a bivariate normal distri-
 744 bution, i.e., $p(\mathbf{s}|\mathbf{u}^f) = \mathcal{N}(\mathbf{s}|\boldsymbol{\mu}_p, \mathbf{K}_p^{-1})$, whose precision matrix \mathbf{K}_p and the mean $\boldsymbol{\mu}_p$ are,

$$\mathbf{K}_p = \boldsymbol{\Lambda}_f + \boldsymbol{\Lambda}_s, \quad \boldsymbol{\mu}_p = \mathbf{K}_p^{-1} \boldsymbol{\Lambda}_f \boldsymbol{\mu}_f. \quad (34)$$

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

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

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

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

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

$$\boldsymbol{\lambda}_{1t} = \mathbf{u}_1^f + \mathbf{u}_{12,t}^r = \mathbf{u}_1^f + w_{12} \mathbf{r}_{2,t-\Delta t}, \quad (36)$$

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

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

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

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

$$\begin{aligned} \mathbf{u}_{12,tj}^r &= U_{12} \exp[-(\theta_j - \tilde{s}_{2,t-\Delta t})^2 / 2a^2] + \delta_\perp \mathbf{u}_{12,tj}^r, \\ \tilde{s}_{2,t-\Delta t} &= \sum_j \mathbf{u}_{12,tj}^r \theta_j / \sum_j \mathbf{u}_{12,tj}^r, \quad \Lambda_s = a^{-2} \sum_j \mathbf{u}_{12,tj}^r, \end{aligned} \quad (39)$$

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

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

772 recurrent weight between two networks is then

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

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

$$\tilde{s}_{mt} = \sum_j \theta_j \mathbf{r}_{mt,j} / \sum_j \mathbf{r}_{mt,j}. \quad (41)$$

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

782 Inference from an information-theoretic point of view

783 The goal of the sampling algorithm is to approximate the posterior distribution of a latent variables,
784 Θ , given a feedforward input, \mathbf{u}^f . Specifically, the latent variables $\Theta = \{s, z\}$ in the hierarchical
785 generative model (Eq. 23), or $\Theta = \mathbf{s} = \{s_1, s_2\}$ in the generative model with breadth (Eq. 33).
786 When the sampling algorithm uses an internal model which does not match the structure of the
787 generative model, the sampling distribution $q(\Theta|\mathbf{u}^f)$ will differ from the true posterior, $p(\Theta|\mathbf{u}^f)$
788 (Eq. 24). In this case the mutual information between the sampling distribution of the latent
789 variables, Θ , and \mathbf{u}^f will be smaller than in the case when samples come from the true posterior,
790 $p(\Theta|\mathbf{u}^f)$,

$$\begin{aligned} I(\Theta, \mathbf{u}^f) &= -\mathbb{E}_{p(\Theta)}[\log p(\Theta)] + \mathbb{E}_{p(\Theta, \mathbf{u}^f)}[\log p(\Theta|\mathbf{u}^f)] \\ &\geq -\mathbb{E}_{p(\Theta)}[\log p(\Theta)] + \mathbb{E}_{p(\Theta, \mathbf{u}^f)}[\log q(\Theta|\mathbf{u}^f)] \equiv I_q(\Theta; \mathbf{u}^f), \end{aligned} \quad (42)$$

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

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

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

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

799 **The neuronal response distribution conditioned on external stimulus**

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

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

806 For simplicity, we only compute the covariability of $p(\mathbf{r}|\mathbf{u}^f)$ along the stimulus subspace (Fig. 1B,
 807 x-axis), because the covariability along other directions is not related with stimulus sampling. By
 808 approximating the Poissonian spiking variability $p(\mathbf{r}|\boldsymbol{\lambda})$ with a multivariate normal distribution
 809 (Eq. 11), and considering the limit of weak fluctuations in $\boldsymbol{\lambda}$ along the stimulus subspace over time,
 810 $p(\mathbf{r}|\mathbf{u}^f)$ can be computed approximately as (see math details in Supplemental Information),

$$\begin{aligned} p(\mathbf{r}|\mathbf{u}^f) &= \int p(\mathbf{r}|\boldsymbol{\lambda})p(\boldsymbol{\lambda}|\mathbf{u}^f)d\boldsymbol{\lambda}, \\ &\approx \mathcal{N}[\mathbf{r}|\mathbf{f}(s), \text{diag}(\mathbf{f}(s)) + V(\bar{s}|\mu_f)\mathbf{f}'_s\mathbf{f}'_s{}^\top], \quad \text{where } s = \mu_f. \end{aligned} \quad (44)$$

811 $\mathbf{f}(s) = \langle \boldsymbol{\lambda}_t \rangle$ denotes the temporally averaged population response. The covariance structure of the
 812 neuronal response includes two terms: $\text{diag}(\mathbf{f}(s))$, a diagonal matrix whose entries equal that of
 813 the vector $\mathbf{f}(s)$ denoting the (independent) Poisson spiking variability (Eq. 23), and $V(\bar{s}|\mu_f)\mathbf{f}'_s\mathbf{f}'_s{}^\top$,
 814 a term that captures the covariability due to firing rate fluctuations along the stimulus subspace
 815 (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
 816 $\mathbf{f}'_s\mathbf{f}'_s{}^\top$ is often termed differential (noise) correlations [4, 17]. With the Gaussian profile of $\mathbf{f}(s)$
 817 (Eqs. 18 and 29), $\mathbf{f}'_s\mathbf{f}'_s{}^\top$ exhibits anti-symmetric structure (Fig. 8B) [17, 22, 50, 68, 69].

818 In Eq. (44), $V(\bar{s}|\mu_f)$ is the variance of \bar{s}_t (the mean of conditional distribution in Eq. 4a) over
 819 time and characterizes the amplitude of internally generated differential correlations. In network
 820 implementation, \bar{s}_t and μ_f are represented as the position of $\boldsymbol{\lambda}_t$ and \mathbf{u}^f on the stimulus subspace
 821 respectively (Eqs. 14 and 20). The dynamics of Gibbs sampling (Eq. S20 in Supplemental Infor-
 822 mation) and the network structure (Eq. 6) imply that

$$V(\bar{s}|\mu_f) = \frac{\Lambda_s}{\Lambda_f(\Lambda_f + \Lambda_s)} = a^2 n_f^{-1} w_E^*. \quad (45)$$

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

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

$$\begin{aligned} p(\mathbf{r}|s) &\approx \mathcal{N}[\mathbf{r}|\mathbf{f}(s), \text{diag}(\mathbf{f}(s)) + V(\bar{s}|s)\mathbf{f}'_s\mathbf{f}'_s{}^\top], \\ V(\bar{s}|s) &= V(\bar{s}|\mu_f) + V(\mu_f|s) = \frac{\Lambda_s}{\Lambda_f(\Lambda_f + \Lambda_s)} + \frac{1}{\Lambda_f} = a^2 n_f^{-1} (w_E^* + 1). \end{aligned} \quad (46)$$

828 Here the variance, $V(\bar{s}|s)$, in the stimulus feature subspace is a mixture of internal variability,
 829 $V(\bar{s}|\mu_f)$, and sensory noise, $V(\mu_f|s)$ (Eq. 23). The neuronal response distribution in coupled net-
 830 works (Fig. 6A) can be obtained similarly (see the Supplemental Information).

831 **A spiking network model with excitatory and inhibitory Poisson neurons**

832 To test the proposed inference mechanisms in a network consisting of E neurons (Eqs. 10-37), we
 833 simulated a well studied recurrently coupled cortical model [21, 22]. The network consisted of N_E
 834 excitatory (E) and N_I inhibitory (I) spiking neurons, with the activity of each neuron modeled as
 835 a Hawkes process [65]. At time t , we represent the response of neuron j in population $a = \{E, I\}$,
 836 \mathbf{r}_{tj}^a , as a spike count drawn from a Poisson distribution with instantaneous firing rate, $\boldsymbol{\lambda}_{tj}^a$,

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

837 Each neuron has a refractory period of 2ms after emitting a spike. The firing rate $\boldsymbol{\lambda}_{tj}^a$ is the sum
 838 of feedforward input \mathbf{u}_{tj}^{af} and recurrent input \mathbf{u}_{tj}^{ar} , so that $\boldsymbol{\lambda}_{tj}^a = \mathbf{u}_{tj}^{af} + \mathbf{u}_{tj}^{ar}$. The feedforward inputs
 839 are filtered spikes from upstream neurons, $\mathbf{u}_{tj}^{af} = \sum_n \eta(t - t_{jn}^f)$, where t_{jn}^f is the time of the n^{th}
 840 spike received by neuron j of population a from the feedforward inputs. Here $\eta(t)$ is the synaptic
 841 input profile which is modeled as $\eta(t) = \exp(-t/\tau_d)/\tau_d$, ($t > 0$). Throughout, we set the synaptic
 842 time constant $\tau_d = 2\text{ms}$. To mimic the Poisson-like variability to sample a context in a hierarchical
 843 generative model (Eqs. 23 and 31), the recurrent input received by neuron j in population a is
 844 defined by

$$\begin{aligned} \mathbf{u}_{tj}^{ar} &= \bar{\mathbf{u}}_{tj}^{ar} + \sqrt{[\bar{\mathbf{u}}_{tj}^{ar}]_+} \boldsymbol{\xi}_t, \\ \bar{\mathbf{u}}_{tj}^{ar} &= \sum_{b=\{E,I\}} \sum_{k=1}^{N_b} \frac{J_{jk}^{ab}}{\sqrt{N}} \sum_n \eta(t - t_{kn}^b), \end{aligned} \quad (48)$$

845 where $\bar{\mathbf{u}}_{tj}^{ar}$ is the mean recurrent input at time t given the neuronal activities of the presynaptic
 846 neurons. The recurrent input in the network is corrupted by noise whose variance equals the mean

847 of the recurrent input. In a physiological network, recurrent noise may be generated by the chaotic
 848 state in network dynamics [70] or synaptic noise [66, 67]. In Eq. (48) the function $[\cdot]_+$ rectifies
 849 the negative input, and ξ_t is a random variable following a standard Gaussian distribution. The
 850 coefficient J_{ij}^{ab} is the synaptic weight from neuron j in population b to neuron i in population a . The
 851 time t_{kn}^b is the time of the n^{th} spike fired by neuron k in population b . The parameter $N = N_E + N_I$
 852 is the total number of neurons in the network. The scaling of the synaptic weights by $1/\sqrt{N}$ is
 853 standard in networks where excitation is balanced by recurrent inhibition [70]. Finally, the synaptic
 854 input profile of the recurrent input, $\eta(t)$, is the same as the one we chose for the feedforward input
 855 for convenience. Note that the rectification in Eq. (48) on recurrent inputs will introduce errors
 856 resulting in deviations of the sampling distribution from the true posterior, and hence we chose
 857 the recurrent weights to be small (Fig. 5). The rectification only arises when using (continuous)
 858 recurrent inputs to sample the context variable, and doesn't impact the generality of sampling by
 859 (discrete) Poisson spiking variability.

860 To model the coding of a circular stimulus such as orientation, the excitatory neurons are
 861 arranged on a ring [22, 68]. The preferred stimuli, θ_j , of the excitatory neurons are equally spaced
 862 on the interval $(-180^\circ, 180^\circ)$, consistent with the range of latent features (Eq. 21). Inhibitory
 863 neurons are not tuned to stimulus, and their role is to stabilize network responses. Note that the
 864 recurrent connections between E neurons are modeled using a Gaussian function decaying with
 865 the distance between the stimuli preferred by the two cells, rather than only self-connection in the
 866 simple network with only E neurons (Eqs. 30),

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

867 We imposed periodic boundaries on the Gaussian function to avoid boundary effect in simulations.
 868 Although in the generative model we assumed non-periodic feature variables (Eq. 3), as long as
 869 the variance of the associated distributions are smaller than the width of the feature space, the
 870 network model with periodic boundaries on the recurrent connection (Eq. 49) provides a good
 871 approximation of the non-periodic Gaussian posterior (Eq. 24). The weight w_{EE} denotes the
 872 average connection strength of all E to E connections. The parameter $a = 40^\circ$ defines the footprint
 873 of connectivity in feature space (i.e the ring), and $L = 360^\circ$ is the length of the ring manifold
 874 (Eq. 21); Multiplication by L in Eq. (49) sets the sum of all E to E connection strengths equal
 875 to $N_E w_{EE}$. Moreover, the excitatory and inhibitory neurons are all-to-all connected with each
 876 other (similar for I to I connections). For simplicity, we consider the E to I , I to I and I to E
 877 connections all to be unstructured (in feature space) and assume that connections of the same type
 878 have equal weight, i.e., $J_{jk}^{EI} = w_{EI}$, $J_{jk}^{IE} = w_{IE}$ and $J_{jk}^{II} = w_{II}$. To simplify the network further,
 879 we consider the connections from the same population of neurons to have the same average weight,
 880 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,

881 we only remove the E recurrent connections between E neurons, i.e., $w_{EE} = 0$, while keeping other
882 connections, including w_{EI} , w_{II} , and w_{IE} , the same as the recurrent network.

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

$$\langle \mathbf{u}_j^{If} \rangle = \frac{w_{If}}{N_I} \sum_{j=1}^{N_E} \langle \mathbf{u}_j^{Ef}(s) \rangle. \quad (50)$$

888 In the simulations, we started with a network of $N_E = 180$ excitatory and $N_I = 45$ inhibitory
889 neurons, and increased the number of neurons by a fixed factor in Fig. 1D. The ratio between the
890 average connection from I neurons and the one from E neurons was kept fixed with $w_I/w_E = 5$.
891 We set the feedforward weight of input to I neurons to $w_{If} = 0.8$. We simulated the dynamics
892 of the model network using the Euler method with a time step of 0.1ms. The typical parameters
893 used in simulation can be found in Table 1 in Supplemental Information. Further details about the
894 simulations and numerical estimates of mutual information and linear Fisher information are also
895 presented in Supplemental Information. The code of network simulation was written in MATLAB
896 2018b, and can be found at GitHub (https://github.com/wenhao-z/Sampling_PoissSpk_Neuron).

897 **A spiking network model of coupled neural circuits**

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

912 Comparing the sampling distribution with posterior in coupled neural circuits

913 We read out the samples from the posterior distribution of each stimulus, \tilde{s}_{mt} , individually from
914 the spiking activities of E neurons, \mathbf{r}_{mt} , in network m in every time window of 20ms by using a
915 population vector. We used this collection of samples to estimate the mean, $\langle \tilde{\mathbf{s}} \rangle = (\langle \tilde{s}_1 \rangle, \langle \tilde{s}_2 \rangle)^\top$,
916 and covariance matrix, Σ_s , of the sampling distribution. Meanwhile, the mean $\boldsymbol{\mu}_f$ and precision
917 matrix Λ_f of the likelihood are linearly read out from the feedforward inputs fed into the network
918 model (Eq. 33).

919 If the sampling distribution is comparable with the posterior, the sampling mean $\langle \tilde{\mathbf{s}} \rangle$ and co-
920 variance Σ_s should satisfy Eq. (34). We use the actual sampling covariance and the likelihood
921 parameters to predict the sampling mean, i.e., $\langle \tilde{\mathbf{s}} \rangle_{\text{pred}} = \Sigma_s \Lambda_f \boldsymbol{\mu}_f$, and compare it with the ac-
922 tual $\langle \tilde{\mathbf{s}} \rangle$ (Fig. 7D-F). To obtain the posterior precision matrix, given the sampling mean $\langle \tilde{\mathbf{s}} \rangle$ and
923 the likelihood parameters, we vary the single parameter of prior precision Λ_s to minimize the KL
924 divergence from the prediction of posterior by using the value of Λ_s , and the actual sampling distri-
925 bution. Given this value of Λ_s , the prediction of posterior precision is computed as $\mathbf{K}_{\text{pred}} = \Lambda_s + \Lambda_f$
926 (Eq. 34) which is then compared with actual sampling precision matrix (Σ_s^{-1} ; see Fig. 7C-G). The
927 prior precision, Λ_s , is a *subjective* prior, which reflects the prior stored in the recurrent network
928 and may change with input (see Discussion). More details of network simulation and parameters
929 can be found in Supplemental Information.

930 Acknowledgements

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

934 References

- 935 [1] Alexandre Pouget, Peter Dayan, and Richard S Zemel. Inference and computation with pop-
936 ulation codes. *Annual Review of Neuroscience*, 26(1):381–410, 2003.
- 937 [2] Brent Doiron, Ashok Litwin-Kumar, Robert Rosenbaum, Gabriel K Ocker, and Krešimir Josić.
938 The mechanics of state-dependent neural correlations. *Nature neuroscience*, 19(3):383–393,
939 2016.
- 940 [3] Robbe LT Goris, J Anthony Movshon, and Eero P Simoncelli. Partitioning neuronal variability.
941 *Nature neuroscience*, 17(6):858–865, 2014.
- 942 [4] Adam Kohn, Ruben Coen-Cagli, Ingmar Kanitscheider, and Alexandre Pouget. Correlations
943 and neuronal population information. *Annual review of neuroscience*, 39:237–256, 2016.
- 944 [5] Julie A Harris, Stefan Mihalas, Karla E Hirokawa, Jennifer D Whitesell, Hannah Choi, Amy
945 Bernard, Phillip Bohn, Shiella Caldejon, Linzy Casal, Andrew Cho, et al. Hierarchical orga-
946 nization of cortical and thalamic connectivity. *Nature*, 575(7781):195–202, 2019.
- 947 [6] Seung Wook Oh, Julie A Harris, Lydia Ng, Brent Winslow, Nicholas Cain, Stefan Mihalas,
948 Quanxin Wang, Chris Lau, Leonard Kuan, Alex M Henry, et al. A mesoscale connectome of
949 the mouse brain. *Nature*, 508(7495):207–214, 2014.
- 950 [7] Rodney J Douglas and Kevan AC Martin. Neuronal circuits of the neocortex. *Annu. Rev.*
951 *Neurosci.*, 27:419–451, 2004.
- 952 [8] L Federico Rossi, Kenneth D Harris, and Matteo Carandini. Spatial connectivity matches
953 direction selectivity in visual cortex. *Nature*, 588(7839):648–652, 2020.
- 954 [9] Kenneth D Harris and Thomas D Mrsic-Flogel. Cortical connectivity and sensory coding.
955 *Nature*, 503(7474):51–58, 2013.
- 956 [10] Marc O Ernst and Martin S Banks. Humans integrate visual and haptic information in a
957 statistically optimal fashion. *Nature*, 415(6870):429–433, 2002.
- 958 [11] Alan Yuille and Daniel Kersten. Vision as bayesian inference: analysis by synthesis? *Trends*
959 *in cognitive sciences*, 10(7):301–308, 2006.
- 960 [12] Tai Sing Lee and David Mumford. Hierarchical bayesian inference in the visual cortex. *JOSA*
961 *A*, 20(7):1434–1448, 2003.
- 962 [13] Konrad P Körding and Daniel M Wolpert. Bayesian integration in sensorimotor learning.
963 *Nature*, 427(6971):244–247, 2004.

- 964 [14] David C Knill and Alexandre Pouget. The bayesian brain: the role of uncertainty in neural
965 coding and computation. *TRENDS in Neurosciences*, 27(12):712–719, 2004.
- 966 [15] Alexandre Pouget, Jeffrey M Beck, Wei Ji Ma, and Peter E Latham. Probabilistic brains:
967 knowns and unknowns. *Nature neuroscience*, 16(9):1170, 2013.
- 968 [16] József Fiser, Pietro Berkes, Gergő Orbán, and Máté Lengyel. Statistically optimal perception
969 and learning: from behavior to neural representations. *Trends in cognitive sciences*, 14(3):119–
970 130, 2010.
- 971 [17] Rubén Moreno-Bote, Jeffrey Beck, Ingmar Kanitscheider, Xaq Pitkow, Peter Latham, and
972 Alexandre Pouget. Information-limiting correlations. *Nature neuroscience*, 17(10):1410, 2014.
- 973 [18] Peter Dayan and Laurence F Abbott. *Theoretical neuroscience*, volume 806. Cambridge, MA:
974 MIT Press, 2001.
- 975 [19] Bruno B Averbeck, Peter E Latham, and Alexandre Pouget. Neural correlations, population
976 coding and computation. *Nature reviews neuroscience*, 7(5):358, 2006.
- 977 [20] Apostolos P Georgopoulos, Andrew B Schwartz, and Ronald E Kettner. Neuronal population
978 coding of movement direction. *Science*, 233(4771):1416–1419, 1986.
- 979 [21] Daniel B Rubin, Stephen D Van Hooser, and Kenneth D Miller. The stabilized supralinear
980 network: a unifying circuit motif underlying multi-input integration in sensory cortex. *Neuron*,
981 85(2):402–417, 2015.
- 982 [22] R Ben-Yishai, R Lev Bar-Or, and H Sompolinsky. Theory of orientation tuning in visual
983 cortex. *Proceedings of the National Academy of Sciences*, 92(9):3844–3848, 1995.
- 984 [23] David C Somers, Sacha B Nelson, and Mriganka Sur. An emergent model of orientation
985 selectivity in cat visual cortical simple cells. *Journal of Neuroscience*, 15(8):5448–5465, 1995.
- 986 [24] Chengcheng Huang, Alexandre Pouget, and Brent David Doiron. Internally generated popu-
987 lation activity in cortical networks hinders information transmission. *bioRxiv*, 2020.
- 988 [25] Daniel Kersten, Pascal Mamassian, and Alan Yuille. Object perception as bayesian inference.
989 *Annu. Rev. Psychol.*, 55:271–304, 2004.
- 990 [26] Kenji Doya, Shin Ishii, Alexandre Pouget, and Rajesh PN Rao. *Bayesian brain: Probabilistic*
991 *approaches to neural coding*. MIT press, 2007.
- 992 [27] Patrik O Hoyer and Aapo Hyvärinen. Interpreting neural response variability as monte carlo
993 sampling of the posterior. In *Advances in neural information processing systems*, pages 293–
994 300, 2003.

- 995 [28] Lars Buesing, Johannes Bill, Bernhard Nessler, and Wolfgang Maass. Neural dynamics as
996 sampling: a model for stochastic computation in recurrent networks of spiking neurons. *PLoS*
997 *computational biology*, 7(11):e1002211, 2011.
- 998 [29] Cristina Savin and Sophie Deneve. Spatio-temporal representations of uncertainty in spiking
999 neural networks. In *NIPS*, volume 27, pages 2024–2032, 2014.
- 1000 [30] Gergő Orbán, Pietro Berkes, József Fiser, and Máté Lengyel. Neural variability and sampling-
1001 based probabilistic representations in the visual cortex. *Neuron*, 92(2):530–543, 2016.
- 1002 [31] Ralf M Haefner, Pietro Berkes, and József Fiser. Perceptual decision-making as probabilistic
1003 inference by neural sampling. *Neuron*, 90(3):649–660, 2016.
- 1004 [32] Rodrigo Echeveste, Laurence Aitchison, Guillaume Hennequin, and Máté Lengyel. Cortical-like
1005 dynamics in recurrent circuits optimized for sampling-based probabilistic inference. *bioRxiv*,
1006 page 696088, 2020.
- 1007 [33] Laurence Aitchison and Máté Lengyel. The hamiltonian brain: Efficient probabilistic in-
1008 ference with excitatory-inhibitory neural circuit dynamics. *PLoS computational biology*,
1009 12(12):e1005186, 2016.
- 1010 [34] Kenneth H Britten, Michael N Shadlen, William T Newsome, and J Anthony Movshon. The
1011 analysis of visual motion: a comparison of neuronal and psychophysical performance. *Journal*
1012 *of Neuroscience*, 12(12):4745–4765, 1992.
- 1013 [35] Christopher M Bishop. *Pattern recognition and machine learning*. springer, 2006.
- 1014 [36] Wei Ji Ma, Jeffrey M Beck, Peter E Latham, and Alexandre Pouget. Bayesian inference with
1015 probabilistic population codes. *Nature Neuroscience*, 9(11):1432–1438, 2006.
- 1016 [37] Mehrdad Jazayeri and J Anthony Movshon. Optimal representation of sensory information by
1017 neural populations. *Nature Neuroscience*, 9(5):690–696, 2006.
- 1018 [38] Bruno A Olshausen and David J Field. Emergence of simple-cell receptive field properties by
1019 learning a sparse code for natural images. *Nature*, 381(6583):607, 1996.
- 1020 [39] Eric R Kandel, James H Schwartz, Thomas M Jessell, Department of Biochemistry, Molecular
1021 Biophysics Thomas Jessell, Steven Siegelbaum, and AJ Hudspeth. *Principles of neural science*,
1022 volume 4. McGraw-hill New York, 2000.
- 1023 [40] Michael S Lewicki and Terrence J Sejnowski. Bayesian unsupervised learning of higher order
1024 structure. *Advances in neural information processing systems*, pages 529–535, 1997.

- 1025 [41] Agnieszka Grabska-Barwinska, Jeffrey M Beck, Alexandre Pouget, and Peter E Latham.
1026 Demixing odorsfast inference in olfaction. *Advances in Neural Information Processing Sys-*
1027 *tems 26 (NIPS 2013)*, 2013.
- 1028 [42] David J Field, Anthony Hayes, and Robert F Hess. Contour integration by the human visual
1029 system: evidence for a local association field. *Vision research*, 33(2):173–193, 1993.
- 1030 [43] Wilson S Geisler, Jeffrey S Perry, BJ Super, and DP Gallogly. Edge co-occurrence in natural
1031 images predicts contour grouping performance. *Vision research*, 41(6):711–724, 2001.
- 1032 [44] Lee Cossell, Maria Florencia Iacaruso, Dylan R Muir, Rachael Houlton, Elie N Sader, Ho Ko,
1033 Sonja B Hofer, and Thomas D Mrsic-Flogel. Functional organization of excitatory synaptic
1034 strength in primary visual cortex. *Nature*, 518(7539):399–403, 2015.
- 1035 [45] Ingmar Kanitscheider, Ruben Coen-Cagli, Adam Kohn, and Alexandre Pouget. Measur-
1036 ing fisher information accurately in correlated neural populations. *PLoS Comput Biol*,
1037 11(6):e1004218, 2015.
- 1038 [46] Tai Sing Lee. The visual system’s internal model of the world. *Proceedings of the IEEE*,
1039 103(8):1359–1378, 2015.
- 1040 [47] Rajkumar Vasudeva Raju and Zachary Pitkow. Inference by reparameterization in neural
1041 population codes. *Advances in Neural Information Processing Systems*, 29:2029–2037, 2016.
- 1042 [48] Sabyasachi Shivkumar, Richard Lange, Ankani Chattoraj, and Ralf Haefner. A probabilistic
1043 population code based on neural samples. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman,
1044 N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*,
1045 volume 31. Curran Associates, Inc., 2018.
- 1046 [49] Ingmar Kanitscheider, Ruben Coen-Cagli, and Alexandre Pouget. Origin of information-
1047 limiting noise correlations. *Proceedings of the National Academy of Sciences*, 112(50):E6973–
1048 E6982, 2015.
- 1049 [50] Adrián Ponce-Alvarez, Alexander Thiele, Thomas D Albright, Gene R Stoner, and Gustavo
1050 Deco. Stimulus-dependent variability and noise correlations in cortical mt neurons. *Proceedings*
1051 *of the National Academy of Sciences*, 110(32):13162–13167, 2013.
- 1052 [51] Si Wu, KY Michael Wong, CC Alan Fung, Yuanyuan Mi, and Wenhao Zhang. Continuous at-
1053 tractor neural networks: candidate of a canonical model for neural information representation.
1054 *F1000Research*, 5, 2016.

- 1055 [52] Guillaume Hennequin, Yashar Ahmadian, Daniel B Rubin, Máté Lengyel, and Kenneth D
1056 Miller. The dynamical regime of sensory cortex: stable dynamics around a single stimulus-
1057 tuned attractor account for patterns of noise variability. *Neuron*, 98(4):846–860, 2018.
- 1058 [53] Richard D Lange and Ralf M Haefner. Task-induced neural covariability as a signature of
1059 approximate bayesian learning and inference. *bioRxiv*, page 081661, 2020.
- 1060 [54] Kechen Zhang. Representation of spatial orientation by the intrinsic dynamics of the head-
1061 direction cell ensemble: a theory. *The Journal of Neuroscience*, 16(6):2112–2126, 1996.
- 1062 [55] Sophie Deneve, Peter E Latham, and Alexandre Pouget. Reading population codes: a neural
1063 implementation of ideal observers. *Nature Neuroscience*, 2(8):740–745, 1999.
- 1064 [56] Si Wu, Shun-ichi Amari, and Hiroyuki Nakahara. Population coding and decoding in a neural
1065 field: a computational study. *Neural Computation*, 14(5):999–1026, 2002.
- 1066 [57] Misha Tsodyks, Klaus Pawelzik, and Henry Markram. Neural networks with dynamic synapses.
1067 *Neural computation*, 10(4):821–835, 1998.
- 1068 [58] Eric Schulz, Joshua B Tenenbaum, David Duvenaud, Maarten Speekenbrink, and Samuel J
1069 Gershman. Compositional inductive biases in function learning. *Cognitive psychology*, 99:44–
1070 79, 2017.
- 1071 [59] Richard D Lange, Ankani Chattoraj, Jeffrey Beck, Jacob Yates, and Ralf Haefner. A confir-
1072 mation bias in perceptual decision-making due to hierarchical approximate inference. *bioRxiv*,
1073 page 440321, 2021.
- 1074 [60] Ruben Coen-Cagli, Adam Kohn, and Odelia Schwartz. Flexible gating of contextual influences
1075 in natural vision. *Nature neuroscience*, 18(11):1648–1655, 2015.
- 1076 [61] Mark M Churchland, M Yu Byron, John P Cunningham, Leo P Sugrue, Marlene R Cohen,
1077 Greg S Corrado, William T Newsome, Andrew M Clark, Paymon Hosseini, Benjamin B Scott,
1078 et al. Stimulus onset quenches neural variability: a widespread cortical phenomenon. *Nature*
1079 *neuroscience*, 13(3):369–378, 2010.
- 1080 [62] Gaby Maimon and John A Assad. Beyond poisson: increased spike-time regularity across
1081 primate parietal cortex. *Neuron*, 62(3):426–440, 2009.
- 1082 [63] Wenhao Zhang, Tai Sing Lee, Brent Doiron, and Si Wu. Distributed sampling-based bayesian
1083 inference in coupled neural circuits. *bioRxiv*, 2020.
- 1084 [64] Deep Ganguli and Eero P Simoncelli. Implicit encoding of prior probabilities in optimal neural
1085 populations. *Advances in neural information processing systems*, 2010:658, 2010.

- 1086 [65] James Trousdale, Yu Hu, Eric Shea-Brown, and Krešimir Josić. Impact of network structure
1087 and cellular response on spike time correlations. *PLoS computational biology*, 8(3):e1002408,
1088 2012.
- 1089 [66] Dmitri A Rusakov, Leonid P Savtchenko, and Peter E Latham. Noisy synaptic conductance:
1090 Bug or a feature? *Trends in Neurosciences*, 2020.
- 1091 [67] Robert Rosenbaum, Jonathan Rubin, and Brent Doiron. Short term synaptic depression
1092 imposes a frequency dependent filter on synaptic information transfer. *PLoS Comput Biol*,
1093 8(6):e1002557, 2012.
- 1094 [68] Si Wu, Kosuke Hamaguchi, and Shun-ichi Amari. Dynamics and computation of continuous
1095 attractors. *Neural Computation*, 20(4):994–1025, 2008.
- 1096 [69] Klaus Wimmer, Duane Q Nykamp, Christos Constantinidis, and Albert Compte. Bump at-
1097 tractor dynamics in prefrontal cortex explains behavioral precision in spatial working memory.
1098 *Nature neuroscience*, 17(3):431, 2014.
- 1099 [70] C van Vreeswijk and Haim Sompolinsky. Chaotic balanced state in a model of cortical circuits.
1100 *Neural computation*, 10(6):1321–1371, 1998.