

Interrogating theoretical models of neural computation with deep inference  
Sean R. Bittner<sup>1</sup>, Agostina Palmigiano<sup>1</sup>, Alex T. Piet<sup>2,3,4</sup>, Chunyu A. Duan<sup>5</sup>, Carlos D. Brody<sup>2,3,6</sup>,  
Kenneth D. Miller<sup>1</sup>, and John P. Cunningham<sup>7</sup>.

<sup>1</sup>Department of Neuroscience, Columbia University,

<sup>2</sup>Princeton Neuroscience Institute,

<sup>3</sup>Princeton University,

<sup>4</sup>Allen Institute for Brain Science,

<sup>5</sup>Institute of Neuroscience, Chinese Academy of Sciences,

<sup>6</sup>Howard Hughes Medical Institute,

<sup>7</sup>Department of Statistics, Columbia University

## <sup>1</sup> 1 Abstract

<sup>2</sup> A cornerstone of theoretical neuroscience is the circuit model: a system of equations that captures a  
<sup>3</sup> hypothesized neural mechanism. Such models are valuable when they give rise to an experimentally  
<sup>4</sup> observed phenomenon – whether behavioral or a pattern of neural activity – and thus can offer  
<sup>5</sup> insights into neural computation. The operation of these circuits, like all models, critically depends  
<sup>6</sup> on the choices of model parameters. A key step is then to identify the model parameters consistent  
<sup>7</sup> with observed phenomena: to solve the inverse problem. To solve challenging inverse problems  
<sup>8</sup> in neuroscience, statistical inference techniques have been used to infer parameters distributions  
<sup>9</sup> most likely to have produced neural datasets. In this work, we present a novel technique, emergent  
<sup>10</sup> property inference (EPI), that brings the power and versatility of the modern probabilistic modeling  
<sup>11</sup> toolkit to theoretical neuroscience. When theorizing circuit models, scientists predominantly focus  
<sup>12</sup> on reproducing computational properties rather than a particular dataset. Our method uses deep  
<sup>13</sup> neural networks to learn parameter distributions with complex structure that produce specific  
<sup>14</sup> computational properties in circuit models. This methodology is introduced through a motivational  
<sup>15</sup> example inferring conductance parameters in a circuit model of the stomatogastric ganglion. Then,  
<sup>16</sup> with recurrent neural networks of increasing size, we show that EPI allows precise control over the  
<sup>17</sup> behavior of inferred parameters, and that EPI scales better in parameter dimension than alternative  
<sup>18</sup> techniques. In the remainder of this work, we present novel theoretical findings gained through  
<sup>19</sup> the examination of complex parametric structure captured by EPI. In a model of primary visual  
<sup>20</sup> cortex, we discovered how connectivity with multiple inhibitory subtypes shapes variability in the

21 excitatory population. Finally, in a model of superior colliculus, we identified and characterized two  
22 distinct regimes of connectivity that facilitate switching between opposite tasks amidst interleaved  
23 trials, mechanistically characterized each regime using probabilistic tools afforded by EPI, and found  
24 conditions where these circuit models reproduce results from optogenetic silencing experiments.  
25 Beyond its scientific contribution, this work illustrates the variety of analyses possible once deep  
26 learning is harnessed towards solving theoretical inverse problems.

## 27 2 Introduction

28 The fundamental practice of theoretical neuroscience is to use a mathematical model to understand  
29 neural computation, whether that computation enables perception, action, or some intermediate  
30 processing. A neural circuit is systematized with a set of equations – the model – and these  
31 equations are motivated by biophysics, neurophysiology, and other conceptual considerations [1–5].  
32 The function of this system is governed by the choice of model *parameters*, which when configured  
33 in a particular way, give rise to a measurable signature of a computation. The work of analyzing  
34 a model then requires solving the inverse problem: given a computation of interest, how can we  
35 reason about the space, manifold, or distribution of parameters that give rise to it? The inverse  
36 problem is crucial for reasoning about likely parameter values, uniquenesses and degeneracies, and  
37 predictions made by the model [6–8].

38 Ideally, one carefully designs a model and analytically derives how computational properties deter-  
39 mine model parameters. Seminal examples of this gold standard include our field’s understanding  
40 of memory capacity in associative neural networks [9], chaos and autocorrelation timescales in ran-  
41 dom neural networks [10], central pattern generation [11], the paradoxical effect [12], and decision  
42 making [13]. Unfortunately, as circuit models include more biological realism, theory via analytical  
43 derivation becomes intractable. Absent this analysis, statistical inference offers a toolkit by which  
44 to solve the inverse problem by identifying, at least approximately, the distribution of parameters  
45 that produce computations in a biologically realistic model [14–19].

46 Statistical inference, of course, requires quantification of the vague term *computation*. In neu-  
47 roscience, two perspectives are dominant. First, often we directly use an *exemplar dataset*: a  
48 collection of samples that express the computation of interest, this data being gathered either ex-  
49 perimentally in the lab or from a computer simulation. While in some sense the best choice given  
50 its connection to experiment [20], some drawbacks exist: these data are well known to have fea-

51 tures irrelevant to the computation of interest [21–23], confounding inferences made on such data.  
52 Related to this point, use of a conventional dataset encourages conventional data likelihoods or loss  
53 functions, which focus on some global metric like squared error or marginal evidence, rather than  
54 the computation itself.

55 Alternatively, researchers often quantify an *emergent property* (EP): a statistic of data that directly  
56 quantifies the computation of interest, wherein the dataset is implicit. While such a choice may  
57 seem esoteric, it is not: the above “gold standard” examples [9–13] all quantify and focus on  
58 some derived feature of the data, rather than the data drawn from the model. An emergent  
59 property is of course a dataset by another name, but it suggests different approach to solving  
60 the same inverse problem: here we directly specify the desired emergent property – a statistic  
61 of data drawn from the model – and the value we wish that property to have, and we set up  
62 an optimization program to find the distribution of parameters that produce this computation.  
63 This statistical framework is not new: it is intimately connected to the literature on approximate  
64 bayesian computation [24–26], parameter sensitivity analyses [27–30], maximum entropy modeling  
65 [31–33], and approximate bayesian inference [34,35]; we detail these connections in Section 5.1.1.

66 The parameter distributions producing a computation may be thin, curved, bent, or multimodal  
67 along various parameter axes and combinations. It is by capturing and quantifying this complex  
68 structure that EPI offers scientific insight. Traditional approximation families (e.g. mean-field or  
69 mixture of gaussians) are limited in the distributional structure they may learn. To address such  
70 restrictions on expressivity, major advancements in machine learning have enabled the use of deep  
71 probability distributions as flexible approximating families for such complicated distributions [36,37]  
72 (see Section 5.1.2). However, the adaptation of deep probability distributions to the problem of  
73 theoretical circuit analysis requires recent developments in deep learning for constrained optimiza-  
74 tion [38], and architectural choices for efficient and expressive deep generative modeling [39, 40].  
75 We detail our method, which we call emergent property inference (EPI) in Section 3.2.

76 Equipped with this method, we demonstrate the capabilities of EPI and present novel theoretical  
77 findings from its analysis. First, we show EPI’s ability to handle biologically realistic circuit models  
78 using a five-neuron model of the stomatogastric ganglion [41]: a neural circuit whose parametric  
79 degeneracy is closely studied [42]. Then, we show EPI’s scalability to high dimensional parameter  
80 distributions by inferring connectivities of recurrent neural networks (RNNs) that exhibit stable,  
81 yet amplified responses – a hallmark of neural responses throughout the brain [43–45]. In a model of  
82 primary visual cortex [46,47], EPI reveals how the recurrent processing across different neuron-type

83 populations shapes excitatory variability: a finding that we show is analytically intractable. Finally,  
84 we investigated the possible connectivities of superior colliculus that allow execution of different  
85 tasks on interleaved trials [48]. EPI discovered a rich distribution containing two connectivity  
86 regimes with different solution classes. We queried the deep probability distribution learned by  
87 EPI to produce a mechanistic understanding of neural responses in each regime. Intriguingly, the  
88 inferred connectivities of each regime reproduced results from optogenetic inactivation experiments  
89 in markedly different ways. These theoretical insights afforded by EPI illustrate the value of deep  
90 inference for the interrogation of neural circuit models.

## 91 3 Results

### 92 3.1 Motivating emergent property inference of theoretical models

93 Consideration of the typical workflow of theoretical modeling clarifies the need for emergent prop-  
94 erty inference. First, one designs or chooses an existing circuit model that, it is hypothesized,  
95 captures the computation of interest. To ground this process in a well-known example, consider  
96 the stomatogastric ganglion (STG) of crustaceans, a small neural circuit which generates multiple  
97 rhythmic muscle activation patterns for digestion [49]. Despite full knowledge of STG connectivity  
98 and a precise characterization of its rhythmic pattern generation, biophysical models of the STG  
99 have complicated relationships between circuit parameters and computation [15, 42].

100 A subcircuit model of the STG [41] is shown schematically in Figure 1A. The fast population (f1  
101 and f2) represents the subnetwork generating the pyloric rhythm and the slow population (s1 and  
102 s2) represents the subnetwork of the gastric mill rhythm. The two fast neurons mutually inhibit  
103 one another, and spike at a greater frequency than the mutually inhibiting slow neurons. The  
104 hub neuron couples with either the fast or slow population, or both depending on modulatory  
105 conditions. The jagged connections indicate electrical coupling having electrical conductance  $g_{el}$ ,  
106 smooth connections in the diagram are inhibitory synaptic projections having strength  $g_{synA}$  onto  
107 the hub neuron, and  $g_{synB} = 5nS$  for mutual inhibitory connections. Note that the behavior of this  
108 model will be critically dependent on its parameterization – the choices of conductance parameters  
109  $\mathbf{z} = [g_{el}, g_{synA}]$ .

110 Second, once the model is selected, one must specify what the model should produce. In this STG  
111 model, we are concerned with neural spiking frequency, which emerges from the dynamics of the  
112 circuit model (Fig. 1B). An emergent property studied by Gutierrez et al. is the hub neuron firing

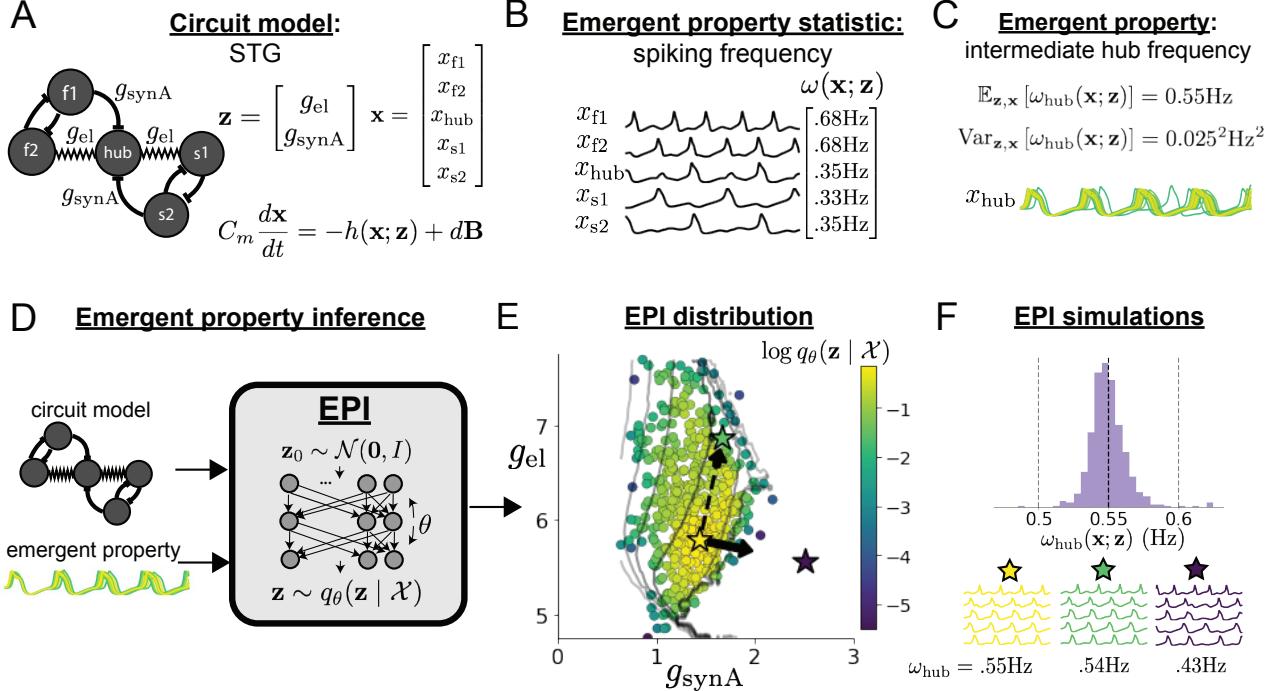


Figure 1: Emergent property inference (EPI) in the stomatogastric ganglion. **A.** Conductance-based subcircuit model of the STG. **B.** Spiking frequency  $\omega(\mathbf{x}; \mathbf{z})$  is an emergent property statistic. Simulated at  $g_{el} = 4.5 \text{nS}$  and  $g_{synA} = 3 \text{nS}$ . **C.** The emergent property of intermediate hub frequency. Simulated activity traces are colored by log probability of generating parameters in the EPI distribution (Panel E). **D.** For a choice of circuit model and emergent property, emergent property inference (EPI) learns a deep probability distribution of parameters  $\mathbf{z}$ . **E.** The EPI distribution producing intermediate hub frequency. Samples are colored by log probability density. Contours of hub neuron frequency error are shown at levels of  $.525, .53, \dots, .575 \text{ Hz}$  (dark to light gray away from mean). Dimension of sensitivity  $\mathbf{v}_1$  (solid arrow) and robustness  $\mathbf{v}_2$  (dashed arrow). **F** (Top) The predictions of the EPI distribution. The black and gray dashed lines show the mean and two standard deviations according the emergent property. (Bottom) Simulations at the starred parameter values.

113 at an intermediate frequency between the intrinsic spiking rates of the fast and slow populations.  
 114 This emergent property is shown in Figure 1C at an average frequency of 0.55Hz. Our notion of  
 115 intermediate hub frequency is not strictly 0.55Hz, but also moderate deviations of this frequency  
 116 between the fast (.35Hz) and slow (.68Hz) frequencies.  
 117 Third, the model parameters producing the emergent property are inferred. By precisely quantify-  
 118 ing the emergent property of interest as a statistical feature of the model, we use EPI to condition  
 119 directly on this emergent property. Before presenting technical details (in the following section), let  
 120 us understand emergent property inference schematically. EPI (Fig. 1D) takes, as input, the model  
 121 and the specified emergent property, and as its output, returns the parameter distribution (Fig.  
 122 1E). This distribution – represented for clarity as samples from the distribution – is a parameter  
 123 distribution constrained such that the circuit model produces the emergent property. Once EPI  
 124 is run, the returned distribution can be used to efficiently generate additional parameter samples.  
 125 Most importantly, the inferred distribution can be efficiently queried to quantify the parametric  
 126 structure that it captures. By quantifying the parametric structure governing the emergent prop-  
 127 erty, EPI informs the central question of this inverse problem: what aspects or combinations of  
 128 model parameters have the desired emergent property?

### 129 3.2 A deep generative modeling approach to emergent property inference

130 Emergent property inference (EPI) formalizes the three-step procedure of the previous section  
 131 with deep probability distributions [36, 37]. First, as is typical, we consider the model as a  
 132 coupled set of noisy differential equations. In this STG example, the model activity (or state)  
 133  $\mathbf{x} = [x_{f1}, x_{f2}, x_{hub}, x_{s1}, x_{s2}]$  is the membrane potential for each neuron, which evolves according to  
 134 the biophysical conductance-based equation:

$$C_m \frac{d\mathbf{x}(t)}{dt} = -h(\mathbf{x}(t); \mathbf{z}) + d\mathbf{B} \quad (1)$$

135 where  $C_m=1nF$ , and  $\mathbf{h}$  is a sum of the leak, calcium, potassium, hyperpolarization, electrical, and  
 136 synaptic currents, all of which have their own complicated dependence on activity  $\mathbf{x}$  and parameters  
 137  $\mathbf{z} = [g_{el}, g_{synA}]$ , and  $d\mathbf{B}$  is white gaussian noise [41] (see Section 5.2.1 for more detail).

138 Second, we determine that our model should produce the emergent property of “intermediate hub  
 139 frequency” (Figure 1C). We stipulate that the hub neuron’s spiking frequency – denoted by statistic  
 140  $\omega_{hub}(\mathbf{x})$  – is close to a frequency of 0.55Hz, between that of the slow and fast frequencies. Mathe-  
 141 matically, we define this emergent property with two constraints: that the mean hub frequency is

142 0.55Hz,

$$\mathbb{E}_{\mathbf{z}, \mathbf{x}} [\omega_{\text{hub}}(\mathbf{x}; \mathbf{z})] = 0.55 \quad (2)$$

143 and that the variance of the hub frequency is moderate

$$\text{Var}_{\mathbf{z}, \mathbf{x}} [\omega_{\text{hub}}(\mathbf{x}; \mathbf{z})] = 0.025^2 \quad (3)$$

144 In the emergent property of intermediate hub frequency, the statistic of hub neuron frequency  
145 is constrained over the distribution of parameters  $\mathbf{z}$  and the distribution of the data  $\mathbf{x}$  that those  
146 parameters produce. This expectation over  $\mathbf{x}$  is taken over any randomness in initial state ( $\mathbf{x}(t=0)$ )  
147 and noise (e.g.  $d\mathbf{B}$  in Equation 1). Formally, the emergent property is the collection of these two  
148 constraints. In general, an emergent property is a collection of constraints on statistical moments  
149 that together define the computation.

150 Third, we perform emergent property inference: we find a distribution over parameter configura-  
151 tions  $\mathbf{z}$  of models that produce the emergent property; in other words, they satisfy the constraints  
152 introduced in Equations 2 and 3. This distribution will be chosen from a family of probability  
153 distributions  $\mathcal{Q} = \{q_{\theta}(\mathbf{z}) : \theta \in \Theta\}$ , defined by a deep neural network [36, 37] (Figure 1D, EPI  
154 box). Deep probability distributions map a simple random variable  $\mathbf{z}_0$  (e.g. an isotropic gaussian)  
155 through a deep neural network with weights and biases  $\theta$  to parameters  $\mathbf{z} = g_{\theta}(\mathbf{z}_0)$  of a suitably  
156 complicated distribution (see Section 5.1.2 for more details). Many distributions in  $\mathcal{Q}$  will respect  
157 the emergent property constraints, so we select the most random (highest entropy) distribution,  
158 which is the same choice commonly made in variational bayesian methods (see Section 5.1.6). In  
159 EPI optimization, stochastic gradient steps in  $\theta$  are taken such that entropy is maximized, and the  
160 emergent property  $\mathcal{X}$  is produced (see Section 5.1). We then denote the inferred EPI distribution  
161 as  $q_{\theta}(\mathbf{z} | \mathcal{X})$ , since the structure of the learned parameter distribution is determined by weights  
162 and biases  $\theta$ , and this distribution is conditioned upon emergent property  $\mathcal{X}$ .

163 The structure of the inferred parameter distributions of EPI can be analyzed to reveal key infor-  
164 mation about how the circuit model produces the emergent property. As probability in the EPI  
165 distribution decreases away from the mode of  $q_{\theta}(\mathbf{z} | \mathcal{X})$  (Fig. 1E yellow star), the emergent prop-  
166 erty deteriorates. Perturbing  $\mathbf{z}$  along a dimension in which  $q_{\theta}(\mathbf{z} | \mathcal{X})$  change little will not disturb  
167 the emergent property, making this parameter combination *robust* with respect to the emergent  
168 property. In contrast, if  $\mathbf{z}$  is perturbed along a dimension with strongly decreasing  $q_{\theta}(\mathbf{z} | \mathcal{X})$ ,  
169 that parameter combination is deemed *sensitive* [27, 30]. By querying the second order derivative  
170 (Hessian) of  $\log q_{\theta}(\mathbf{z} | \mathcal{X})$  at a mode, we can quantitatively identify how sensitive (or robust) each

171 eigenvector is by its eigenvalue; the more negative, the more sensitive and the closer to zero, the  
172 more robust (see Section 5.2.4). Indeed, samples equidistant from the mode along these dimensions  
173 of sensitivity ( $\mathbf{v}_1$ , smaller eigenvalue) and robustness ( $\mathbf{v}_2$ , greater eigenvalue) (Fig. 1E, arrows)  
174 agree with error contours (Fig. 1E contours) and have diminished or preserved hub frequency, re-  
175 spectively (Fig. 1F activity traces). The directionality of  $\mathbf{v}_2$  suggests that changes in conductance  
176 along this parameter combination will most preserve hub neuron firing between the intrinsic rates  
177 of the pyloric and gastric mill rhythms. Importantly, once an EPI distribution has been learned,  
178 the modes and Hessians of the distribution can be measured with trivial computation (see Section  
179 5.1.2).

180 In the following sections, we demonstrate EPI on three neural circuit models across ranges of  
181 biological realism, neural system function, and network scale. First, we demonstrate the superior  
182 scalability of EPI compared to alternative techniques by inferring high-dimensional distributions  
183 of recurrent neural network connectivities that exhibit amplified, yet stable responses. Next, in a  
184 model of primary visual cortex [46,47], we show how EPI discovers parametric degeneracy, revealing  
185 how input variability across neuron types affects the excitatory population. Finally, in a model of  
186 superior colliculus [48], we used EPI to capture multiple parametric regimes of task switching, and  
187 queried the dimensions of parameter sensitivity to characterize each regime.

### 188 3.3 Scaling inference of recurrent neural network connectivity with EPI

189 To understand how EPI scales in comparison to existing techniques, we consider recurrent neu-  
190 ral networks (RNNs). Transient amplification is a hallmark of neural activity throughout cortex,  
191 and is often thought to be intrinsically generated by recurrent connectivity in the responding cor-  
192 tical area [43–45]. It has been shown that to generate such amplified, yet stabilized responses,  
193 the connectivity of RNNs must be non-normal [43,50], and satisfy additional constraints [51]. In  
194 theoretical neuroscience, RNNs are optimized and then examined to show how dynamical systems  
195 could execute a given computation [52,53], but such biologically realistic constraints on connec-  
196 tivity [43,50,51] are ignored for simplicity or because constrained optimization is difficult. In  
197 general, access to distributions of connectivity that produce theoretical criteria like stable amplifi-  
198 cation, chaotic fluctuations [10], or low tangling [54] would add scientific value to existing research  
199 with RNNs. Here, we use EPI to learn RNN connectivities producing stable amplification, and  
200 demonstrate the superior scalability and efficiency of EPI to alternative approaches.

201 We consider a rank-2 RNN with  $N$  neurons having connectivity  $W = UV^\top$  and dynamics

$$\tau \dot{\mathbf{x}} = -\mathbf{x} + W\mathbf{x}, \quad (4)$$

202 where  $U = [\mathbf{U}_1 \ \mathbf{U}_2] + g\chi^{(U)}$ ,  $V = [\mathbf{V}_1 \ \mathbf{V}_2] + g\chi^{(V)}$ ,  $\mathbf{U}_1\mathbf{U}_2, \mathbf{V}_1, \mathbf{V}_2 \in [-1, 1]^N$ , and  $\chi_{i,j}^{(U)}, \chi_{i,j}^{(V)} \sim$   
 203  $\mathcal{N}(0, 1)$ . We infer connectivity parameters  $\mathbf{z} = [\mathbf{U}_1^\top, \mathbf{U}_2^\top, \mathbf{V}_1^\top, \mathbf{V}_2^\top]^\top$  that produce stable amplifi-  
 204 cation. Two conditions are necessary and sufficient for RNNs to exhibit stable amplification [51]:  
 205  $\text{real}(\lambda_1) < 1$  and  $\lambda_1^s > 1$ , where  $\lambda_1$  is the eigenvalue of  $W$  with greatest real part and  $\lambda^s$  is the max-  
 206 imum eigenvalue of  $W^s = \frac{W+W^\top}{2}$ . RNNs with  $\text{real}(\lambda_1) = 0.5 \pm 0.5$  and  $\lambda_1^s = 1.5 \pm 0.5$  will be stable  
 207 with modest decay rate ( $\text{real}(\lambda_1)$  close to its upper bound of 1) and exhibit modest amplification  
 208 ( $\lambda_1^s$  close to its lower bound of 1). EPI can naturally condition on this emergent property

$$\begin{aligned} \mathcal{X} : \mathbb{E}_{\mathbf{z}, \mathbf{x}} \begin{bmatrix} \text{real}(\lambda_1) \\ \lambda_1^s \end{bmatrix} &= \begin{bmatrix} 0.5 \\ 1.5 \end{bmatrix} \\ \text{Var}_{\mathbf{z}, \mathbf{x}} \begin{bmatrix} \text{real}(\lambda_1) \\ \lambda_1^s \end{bmatrix} &= \begin{bmatrix} 0.25^2 \\ 0.25^2 \end{bmatrix}. \end{aligned} \quad (5)$$

209 Variance constraints predicate that the majority of the distribution (within two standard devia-  
 210 tions) are within the specified ranges.

211 For comparison, we infer the parameters  $\mathbf{z}$  likely to produce stable amplification using two al-  
 212 ternative simulation-based inference approaches. Sequential Monte Carlo approximate bayesian  
 213 computation (SMC-ABC) [26] is a rejection sampling approach that uses SMC techniques to im-  
 214 prove efficiency, and sequential neural posterior estimation (SNPE) [35] approximates posteriors  
 215 with deep probability distributions using a two-network architecture (see Section 5.1.1). Unlike  
 216 EPI, these statistical inference techniques do not constrain the statistics of the predictive distribu-  
 217 tion, so they were run by conditioning on an exemplary dataset  $\mathbf{x}_0 = \boldsymbol{\mu}$ , following standard practice  
 218 with these methods [26, 35]. To compare the efficiency of these different techniques, we measured  
 219 the time and number of simulations necessary for the distance of the predictive mean to be less  
 220 than 0.5 from  $\boldsymbol{\mu} = \mathbf{x}_0$  (see Section 5.3).

221 As the number of neurons  $N$  in the RNN, and thus the dimension of the parameter space  $\mathbf{z} \in$   
 222  $[-1, 1]^{4N}$ , is scaled, we see that EPI converges at greater speed and at greater dimension than  
 223 SMC-ABC and SNPE (Fig. 2A). It also becomes most efficient to use EPI in terms of simulation  
 224 count at  $N = 50$  (Fig. 2B). It is well known that ABC techniques struggle in parameter spaces  
 225 of modest dimension [55], yet we were careful to assess the scalability of SNPE, which is a more  
 226 closely related methodology to EPI. Between EPI and SNPE, we closely controlled the number of

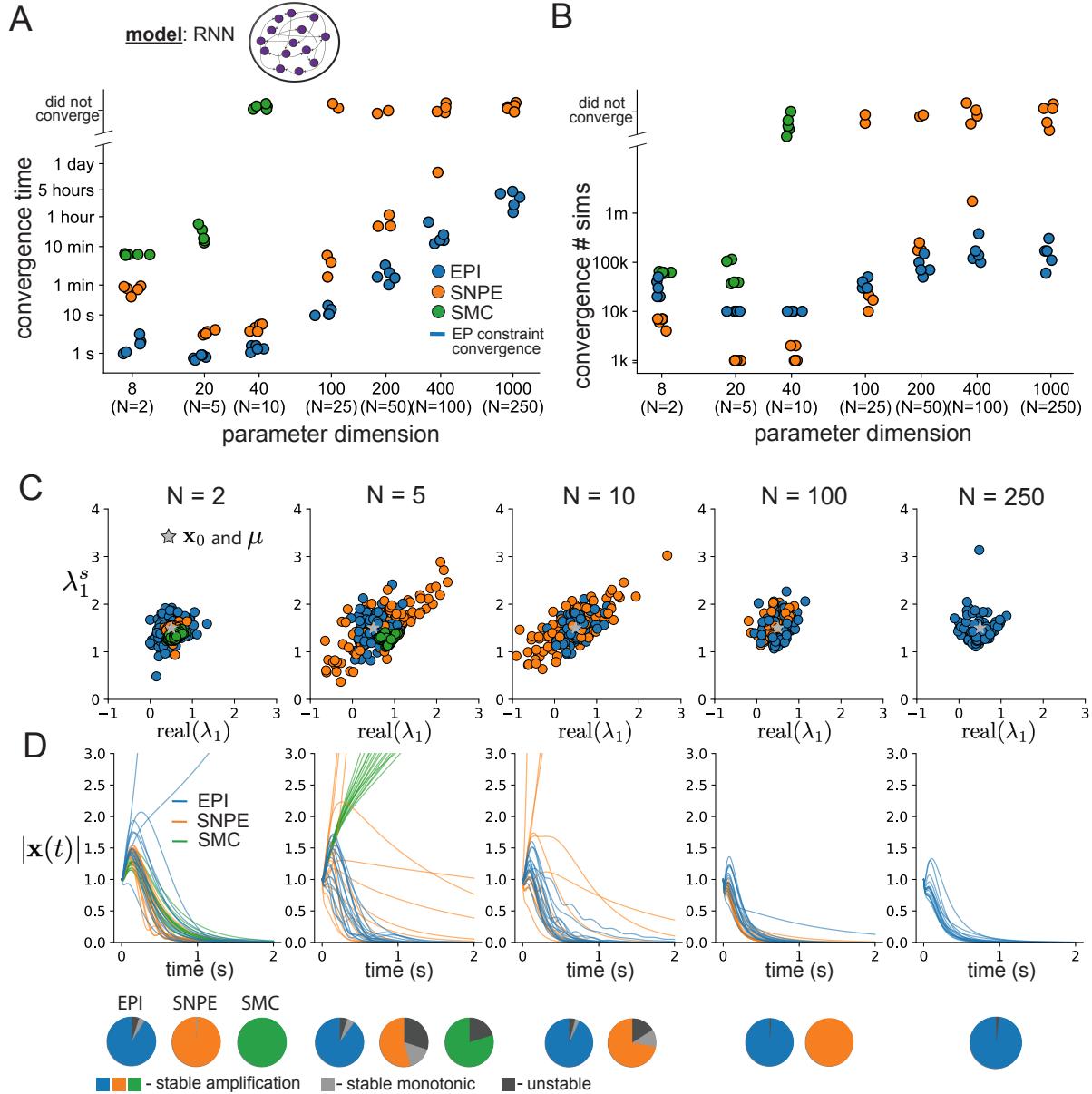


Figure 2: **A.** Wall time of EPI (blue), SNPE (orange), and SMC-ABC (green) to converge on RNN connectivities producing stable amplification. Each dot shows convergence time for an individual random seed. For reference, the mean wall time for EPI to achieve its full constraint convergence (means and variances) is shown (blue line). **B.** Simulation count of each algorithm to achieve convergence. Same conventions as A. **C.** The predictive distributions of connectivities inferred by EPI (blue), SNPE (orange), and SMC-ABC (green), with reference to  $x_0 = \mu$  (gray star). **D.** Simulations of networks inferred by each method ( $\tau = 100ms$ ). Each trace (15 per algorithm) corresponds to simulation of one  $z$ . (Below) Ratio of obtained samples producing stable amplification, monotonic decay, and instability.

parameters in deep probability distributions by dimensionality (Fig. S5), and tested more aggressive SNPE hyperparameter choices when SNPE failed to converge (Fig. S6). In this analysis, we see that deep inference techniques EPI and SNPE are far more amenable to inference of high dimensional RNN connectivities than rejection sampling techniques like SMC-ABC, and that EPI outperforms SNPE in both wall time (elapsed real time) and simulation count.

No matter the number of neurons, EPI always produces connectivity distributions with mean and variance of  $\text{real}(\lambda_1)$  and  $\lambda_1^s$  according to  $\mathcal{X}$  (Fig. 2C, blue). For the dimensionalities in which SMC-ABC is tractable, the inferred parameters are concentrated and offset from the exemplary dataset  $\mathbf{x}_0$  (Fig. 2C, green). When using SNPE, the predictions of the inferred parameters are highly concentrated at some RNN sizes and widely varied in others (Fig. 2C, orange). We see these properties reflected in simulations from the inferred distributions: EPI produces a consistent variety of stable, amplified activity norms  $|\mathbf{x}(t)|$ , SMC-ABC produces a limited variety of responses, and the changing variety of responses from SNPE emphasizes the control of EPI on parameter predictions (Fig. 2D). Even for moderate neuron counts, the predictions of the inferred distribution of SNPE are highly dependent on  $N$  and  $g$ , while EPI maintains the emergent property across choices of RNN (see Section 5.3.5).

To understand these differences, note that EPI outperforms SNPE in high dimensions by using gradient information (from  $\nabla_{\mathbf{z}} f(\mathbf{x}; \mathbf{z}) = \nabla_{\mathbf{z}} [\text{real}(\lambda_1), \lambda_1^s]^\top$ ). This choice agrees with recent speculation that such gradient information could improve the efficiency of simulation-based inference techniques [56], as well as reflecting the classic tradeoff between gradient-based and sampling-based estimators (scaling and speed versus generality). Since gradients of the emergent property statistics are necessary in EPI optimization, gradient tractability is a key criteria when determining the suitability of a simulation-based inference technique. If the emergent property gradient is efficiently calculated, EPI is a clear choice for inferring high dimensional parameter distributions. In the next two sections, we use EPI for novel scientific insight by examining the structure of inferred distributions.

### 3.4 EPI reveals how recurrence with multiple inhibitory subtypes governs excitatory variability in a V1 model

Dynamical models of excitatory (E) and inhibitory (I) populations with supralinear input-output function have succeeded in explaining a host of experimentally documented phenomena in primary visual cortex (V1). In a regime characterized by inhibitory stabilization of strong recurrent excita-

tion, these models give rise to paradoxical responses [12], selective amplification [43, 50], surround suppression [57] and normalization [58]. Recent theoretical work [59] shows that stabilized E-I models reproduce the effect of variability suppression [60]. Furthermore, experimental evidence shows that inhibition is composed of distinct elements – parvalbumin (P), somatostatin (S), VIP (V) – composing 80% of GABAergic interneurons in V1 [61–63], and that these inhibitory cell types follow specific connectivity patterns (Fig. 3A) [64]. Here, we use EPI on a model of V1 with biologically realistic connectivity to show how the structure of input across neuron types affects the variability of the excitatory population – the population largely responsible for projecting to other brain areas [65].

We considered response variability of a nonlinear dynamical V1 circuit model (Fig. 3A) with a state comprised of each neuron-type population’s rate  $\mathbf{x} = [x_E, x_P, x_S, x_V]^\top$ . Each population receives recurrent input  $W\mathbf{x}$ , where  $W$  is the effective connectivity matrix (see Section 5.4) and an external input with mean  $\mathbf{h}$ , which determines population rate via supralinear nonlinearity  $\phi(\cdot) = [\cdot]_+^2$ . The external input has an additive noisy component  $\epsilon$  with variance  $\sigma^2 = [\sigma_E^2, \sigma_P^2, \sigma_S^2, \sigma_V^2]$ . This noise has a slower dynamical timescale  $\tau_{\text{noise}} > \tau$  than the population rate, allowing fluctuations around a stimulus-dependent steady-state (Fig. 3B). This model is the stochastic stabilized supralinear network (SSSN) [59]

$$\tau \frac{d\mathbf{x}}{dt} = -\mathbf{x} + \phi(W\mathbf{x} + \mathbf{h} + \epsilon), \quad (6)$$

generalized to have multiple inhibitory neuron types. It introduces stochasticity to four neuron-type models of V1 [46]. Stochasticity and inhibitory multiplicity introduce substantial complexity to the mathematical treatment of this problem (see Section 5.4.5) motivating the analysis of this model with EPI. Here, we consider fixed weights  $W$  and input  $\mathbf{h}$  [47], and study the effect of input variability  $\mathbf{z} = [\sigma_E, \sigma_P, \sigma_S, \sigma_V]^\top$  on excitatory variability.

We quantify levels of E-population variability by studying two emergent properties

$$\begin{aligned} \mathcal{X}(5 \text{ Hz}) : \mathbb{E}_{\mathbf{z}} [s_E(\mathbf{x}; \mathbf{z})] &= 5 \text{ Hz} & \mathcal{X}(10 \text{ Hz}) : \mathbb{E}_{\mathbf{z}} [s_E(\mathbf{x}; \mathbf{z})] &= 10 \text{ Hz} \\ \text{Var}_{\mathbf{z}} [s_E(\mathbf{x}; \mathbf{z})] &= 1 \text{ Hz}^2 & \text{Var}_{\mathbf{z}} [s_E(\mathbf{x}; \mathbf{z})] &= 1 \text{ Hz}^2, \end{aligned} \quad (7)$$

where  $s_E(\mathbf{x}; \mathbf{z})$  is the standard deviation of the stochastic E-population response about its steady state (Fig. 3C). In the following analyses, we select  $1 \text{ Hz}^2$  variance such that the two emergent properties do not overlap in  $s_E(\mathbf{z}; \mathbf{x})$ .

First, we ran EPI to obtain parameter distribution  $q_{\theta}(\mathbf{z} | \mathcal{X}(5 \text{ Hz}))$  producing E-population variability around 5 Hz (Fig. 3D). From the marginal distribution of  $\sigma_E$  and  $\sigma_P$  (Fig. 3D, top-left),

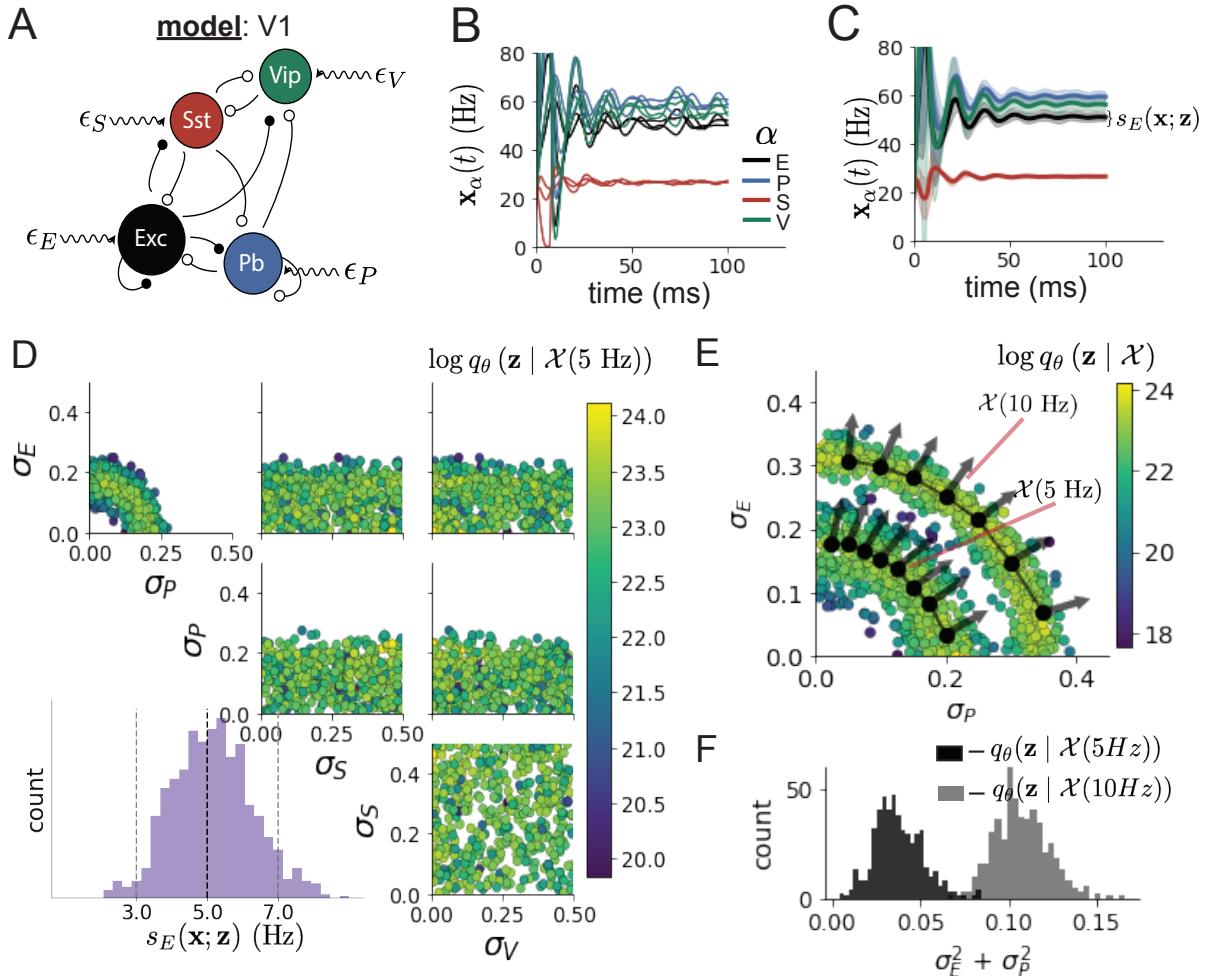


Figure 3: Emergent property inference in the stochastic stabilized supralinear network (SSSN)

**A.** Four-population model of primary visual cortex with excitatory (black), parvalbumin (blue), somatostatin (red), and VIP (green) neurons (excitatory and inhibitory projections filled and unfilled, respectively). Some neuron-types largely do not form synaptic projections to others ( $|W_{\alpha_1, \alpha_2}| < 0.025$ ). Each neural population receives a baseline input  $\mathbf{h}_b$ , and the E- and P- populations also receive a contrast-dependent input  $\mathbf{h}_c$ . Additionally, each neural population receives a slow noisy input  $\epsilon$ . **B.** Transient network responses of the SSSN model. Traces are independent trials with varying initialization  $\mathbf{x}(0)$  and noise  $\epsilon$ . **C.** Mean (solid line) and standard deviation  $s_E(\mathbf{x}; \mathbf{z})$  (shading) across 100 trials. **D.** EPI distribution of noise parameters  $\mathbf{z}$  conditioned on E-population variability. The EPI predictive distribution of  $s_E(\mathbf{x}; \mathbf{z})$  is show on the bottom-left. **E.** (Top) Enlarged visualization of the  $\sigma_E$ - $\sigma_P$  marginal distribution of EPI  $q_\theta(\mathbf{z} | \mathcal{X}(5 \text{ Hz}))$  and  $q_\theta(\mathbf{z} | \mathcal{X}(10 \text{ Hz}))$ . Each black dot shows the mode at each  $\sigma_P$ . The arrows show the most sensitive dimensions of the Hessian evaluated at these modes. **F.** The predictive distributions of  $\sigma_E^2 + \sigma_P^2$  of each inferred distribution  $q_\theta(\mathbf{z} | \mathcal{X}(5 \text{ Hz}))$  and  $q_\theta(\mathbf{z} | \mathcal{X}(10 \text{ Hz}))$ .

286 we can see that  $s_E(\mathbf{x}; \mathbf{z})$  is sensitive to various combinations of  $\sigma_E$  and  $\sigma_P$ . Alternatively, both  $\sigma_S$   
287 and  $\sigma_V$  are degenerate with respect to  $s_E(\mathbf{x}; \mathbf{z})$  evidenced by the, unexpectedly, high variability in  
288 those dimensions (Fig. 3D, bottom-right). Together, these observations imply a curved path with  
289 respect to  $s_E(\mathbf{x}; \mathbf{z})$  of 5 Hz, which is indicated by the modes along  $\sigma_P$  (Fig. 3E).

290 Figure 3E suggests a quadratic relationship in E-population fluctuations and the standard deviation  
291 of E- and P-population input; as the square of either  $\sigma_E$  or  $\sigma_P$  increases, the other compensates by  
292 decreasing to preserve the level of  $s_E(\mathbf{x}; \mathbf{z})$ . This quadratic relationship is preserved at greater level  
293 of E-population variability  $\mathcal{X}(10 \text{ Hz})$  (Fig. 3E and S8). Indeed, the sum of squares of  $\sigma_E$  and  $\sigma_P$  is  
294 larger in  $q_{\theta}(\mathbf{z} | \mathcal{X}(10 \text{ Hz}))$  than  $q_{\theta}(\mathbf{z} | \mathcal{X}(5 \text{ Hz}))$  (Fig 3F,  $p < 1 \times 10^{-10}$ ), while the sum of squares of  
295  $\sigma_S$  and  $\sigma_V$  are not significantly different in the two EPI distributions (Fig. S10,  $p = .40$ ), in which  
296 parameters were bounded from 0 to 0.5. The strong interaction between E- and P-population input  
297 variability on excitatory variability is intriguing, since this circuit exhibits a paradoxical effect in  
298 the P-population (and no other inhibitory types) (Fig. S11), meaning that the E-population is  
299 P-stabilized. Future research may uncover a link between the population of network stabilization  
300 and compensatory interactions governing excitatory variability.

301 EPI revealed the quadratic dependence of excitatory variability on input variability to the E- and  
302 P-populations, as well as its independence to input from the other two inhibitory populations. In a  
303 simplified model ( $\tau = \tau_{\text{noise}}$ ), it can be shown that surfaces of equal variance are ellipsoids as a func-  
304 tion of  $\sigma$  (see Section 5.4.5). Nevertheless, the sensitive and degenerate parameters are challenging  
305 to predict mathematically, since the covariance matrix depends on the steady-state solution of the  
306 network [59, 66], and terms in the covariance expression increase quadratically with each additional  
307 neuron-type population (see also Section 5.4.5). By pointing out this mathematical complexity, we  
308 emphasize the value of streamlined methods for gaining understanding about theoretical models  
309 when mathematical analysis becomes onerous or impractical. While we have just shown that EPI  
310 can be used to investigate fundamental aspects of sensory computation, in the next two sections,  
311 we use the probabilistic tools of EPI to identify and characterize two distinct parametric regimes of  
312 a neural circuit executing a computation, and then relate these insights to behavioral experiments.

### 313 3.5 EPI identifies two regimes of rapid task switching

314 It has been shown that rats can learn to switch from one behavioral task to the next on randomly  
315 interleaved trials [67], and an important question is what neural mechanisms produce this compu-  
316 tation. In this experimental setup, rats were given an explicit task cue on each trial, either Pro

317 or Anti. After a delay period, rats were shown a stimulus, and made a context (task) dependent  
 318 response (Fig. 4A). In the Pro task, rats were required to orient towards the stimulus, while in  
 319 the Anti task, rats were required to orient away from the stimulus. Pharmacological inactivation  
 320 of the SC impaired rat performance, and time-specific optogenetic inactivation revealed a crucial  
 321 role for the SC on the cognitively demanding Anti trials [48]. These results motivated a nonlinear  
 322 dynamical model of the SC containing four functionally-defined neuron-type populations. In Duan  
 323 et al. 2019, a computationally intensive procedure was used to obtain a set of 373 connectivity  
 324 parameters that qualitatively reproduced these optogenetic inactivation results. To build upon  
 325 the insights of this previous work, we use the probabilistic tools afforded by EPI to identify and  
 326 mechanistically characterize two linked, yet distinct regimes of rapid task switching connectivity.

327 In this SC model, there are Pro- and Anti-populations in each hemisphere (left (L) and right (R))  
 328 with activity variables  $\mathbf{x} = [x_{LP}, x_{LA}, x_{RP}, x_{RA}]^\top$  [48]. The connectivity of these populations is  
 329 parameterized by self  $sW$ , vertical  $vW$ , diagonal  $dW$  and horizontal  $hW$  connections (Fig. 4B). The  
 330 input  $\mathbf{h}$  is comprised of a positive cue-dependent signal to the Pro or Anti populations, a positive  
 331 stimulus-dependent input to either the Left or Right populations, and a choice-period input to the  
 332 entire network (see Section 5.5.1). Model responses are bounded from 0 to 1 as a function  $\phi$  of an  
 333 internal variable  $\mathbf{u}$

$$\begin{aligned} \tau \frac{d\mathbf{u}}{dt} &= -\mathbf{u} + W\mathbf{x} + \mathbf{h} + d\mathbf{B} \\ \mathbf{x} &= \phi(\mathbf{u}). \end{aligned} \tag{8}$$

334 The model responds to the side with greater Pro neuron activation; e.g. the response is left if  
 335  $x_{LP} > x_{RP}$  at the end of the trial. Here, we use EPI to determine the network connectivity  
 336  $\mathbf{z} = [sW, vW, dW, hW]^\top$  that produces rapid task switching.

337 Rapid task switching is formalized mathematically as an emergent property with two statistics:  
 338 accuracy in the Pro task  $p_P(\mathbf{x}; \mathbf{z})$  and Anti task  $p_A(\mathbf{x}; \mathbf{z})$ . We stipulate that accuracy be on average  
 339 .75 in each task with variance  $.075^2$

$$\begin{aligned} \mathcal{X} : \mathbb{E}_{\mathbf{z}} \begin{bmatrix} p_P(\mathbf{x}; \mathbf{z}) \\ p_A(\mathbf{x}; \mathbf{z}) \end{bmatrix} &= \begin{bmatrix} .75 \\ .75 \end{bmatrix} \\ \text{Var}_{\mathbf{z}} \begin{bmatrix} p_P(\mathbf{x}; \mathbf{z}) \\ p_A(\mathbf{x}; \mathbf{z}) \end{bmatrix} &= \begin{bmatrix} .075^2 \\ .075^2 \end{bmatrix}. \end{aligned} \tag{9}$$

340 75% accuracy is a realistic level of performance in each task, and with the chosen variance, inferred  
 341 models will not exhibit fully random responses (50%), nor perfect performance (100%).

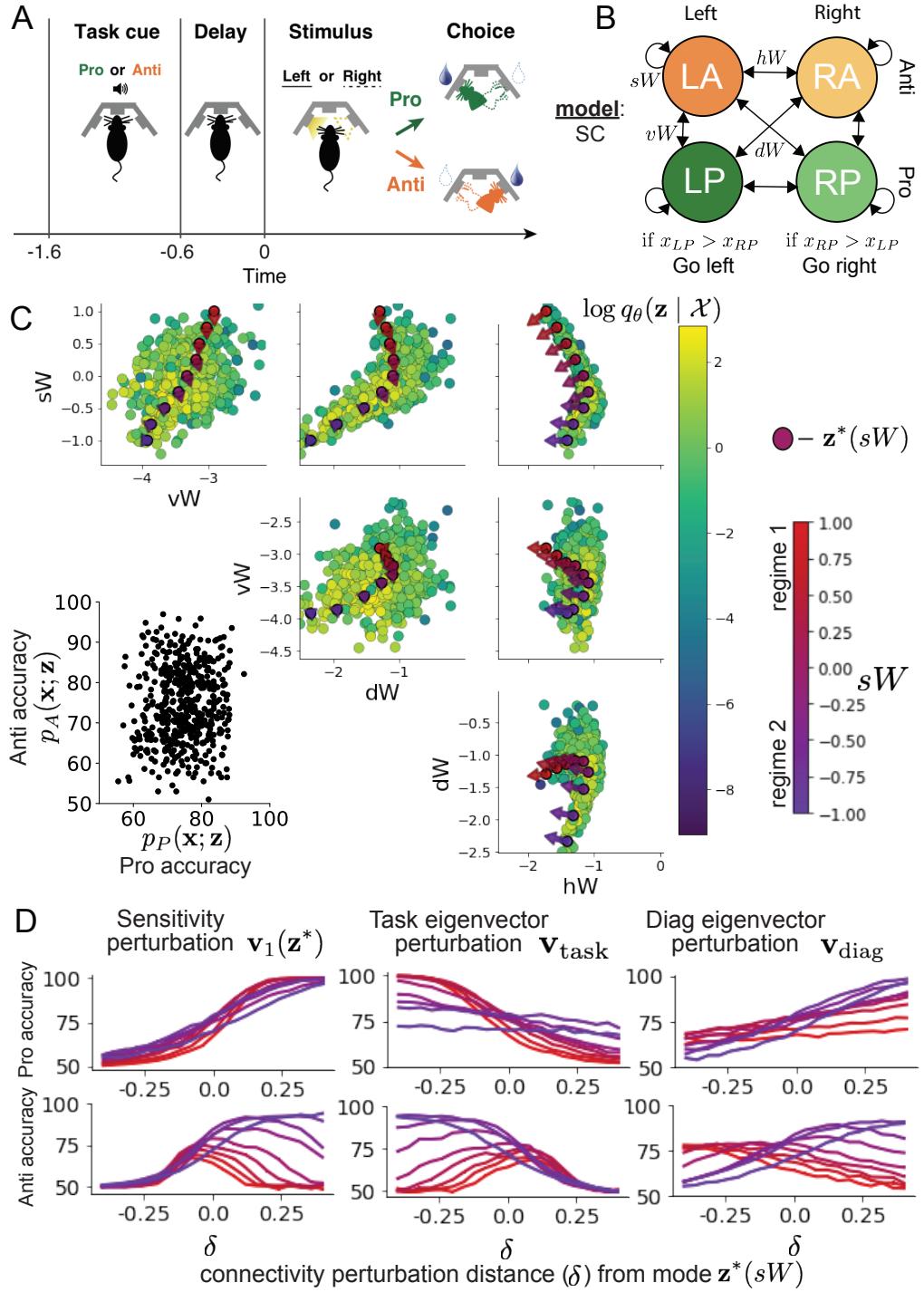


Figure 4: **A.** Rapid task switching behavioral paradigm (see text). **B.** Model of superior colliculus (SC). Neurons: LP - Left Pro, RP - Right Pro, LA - Left Anti, RA - Right Anti. Parameters:  $sW$  - self,  $hW$  - horizontal,  $vW$  - vertical,  $dW$  - diagonal weights. **C.** The EPI inferred distribution of rapid task switching networks. Red/purple parameters indicate modes  $\mathbf{z}^*(sW)$  colored by  $sW$ . Sensitivity vectors  $\mathbf{v}_1(\mathbf{z}^*)$  are shown by arrows. (Bottom-left) EPI predictive distribution of task accuracies. **D.** Mean and standard error ( $N_{\text{test}} = 25$ ) of accuracy in Pro (top) and Anti (bottom) tasks after perturbing connectivity away from mode along  $\mathbf{v}_1(\mathbf{z}^*)$  (left),  $\mathbf{v}_{\text{task}}$  (middle), and  $\mathbf{v}_{\text{diag}}$  (right).

342 The EPI inferred distribution (Fig. 4C) produces Pro and Anti task accuracies (Fig. 4C, bottom-  
 343 left) consistent with rapid task switching (Equation 9). This parameter distribution has rich struc-  
 344 ture, that is not captured well by simple linear correlations (Fig. S12). Specifically, the shape  
 345 of the EPI distribution is sharply bent, matching ground truth structure indicated by brute-force  
 346 sampling (Fig. S18). This is most saliently observed in the marginal distribution of  $sW-hW$  (Fig.  
 347 4C top-right), where anticorrelation between  $sW$  and  $hW$  switches to correlation with decreasing  
 348  $sW$ . By identifying the modes of the EPI distribution  $\mathbf{z}^*(sW)$  at different values of  $sW$  (Fig. 4C  
 349 red/purple dots), we can quantify this change in distributional structure with the sensitivity dimen-  
 350 sion  $\mathbf{v}_1(\mathbf{z})$  (Fig. 4C red/purple arrows). Note that the directionality of these sensitivity dimensions  
 351 at  $\mathbf{z}^*(sW)$  changes distinctly with  $sW$ , and are perpendicular to the flat dimensions of the EPI  
 352 distribution that preserve rapid task switching. These two directionalities of sensitivity motivate  
 353 the distinction of connectivity into two regimes, which produce different types of responses in the  
 354 Pro and Anti tasks (Fig. S13).

355 When perturbing connectivity along the sensitivity dimension away from the modes

$$\mathbf{z} = \mathbf{z}^*(sW) + \delta\mathbf{v}_1(\mathbf{z}^*(sW)), \quad (10)$$

356 Pro accuracy monotonically increases in both regimes (Fig. 4D, top-left). However, there is a stark  
 357 difference between regimes in Anti accuracy. Anti accuracy falls in either direction of  $\mathbf{v}_1$  in regime 1,  
 358 yet monotonically increases along with Pro accuracy in regime 2 (Fig. 4D, bottom-left). The sharp  
 359 change in local structure of the EPI distribution is therefore explained by distinct sensitivities:  
 360 Anti accuracy diminishes in only one or both directions of the sensitivity perturbation.

361 To understand the mechanisms differentiating the two regimes, we can make connectivity pertur-  
 362 bations along dimensions that only modify a single eigenvalue of the connectivity matrix. These  
 363 eigenvalues  $\lambda_{\text{all}}$ ,  $\lambda_{\text{side}}$ ,  $\lambda_{\text{task}}$ , and  $\lambda_{\text{diag}}$  correspond to connectivity eigenmodes with intuitive roles  
 364 in processing in this task (Fig. S14A). For example, greater  $\lambda_{\text{task}}$  will strengthen internal repre-  
 365 sentations of task, while greater  $\lambda_{\text{diag}}$  will amplify dominance of Pro and Anti pairs in opposite  
 366 hemispheres (Section 5.5.7). Unlike the sensitivity dimension, these perturbations of connectivity  
 367 eigenvalue are independent of  $\mathbf{z}^*(sW)$  (see Section 5.5.7), e.g.

$$\mathbf{z} = \mathbf{z}^*(sW) + \delta\mathbf{v}_{\text{task}}. \quad (11)$$

368 Connectivity perturbation analyses reveal that decreasing  $\lambda_{\text{task}}$  has a very similar effect on Anti  
 369 accuracy as perturbations along the sensitivity dimension (Fig. 4D, middle). The similar effects

370 of  $\mathbf{v}_1(\mathbf{z}^*)$  and  $-\mathbf{v}_{\text{task}}$  suggest that there is a carefully tuned strength of task representation in con-  
 371 nectivity regime 1, which if disturbed results in random Anti trial responses. Finally, we recognize  
 372 that increasing  $\lambda_{\text{diag}}$  has opposite effects on Anti accuracy in each regime (Fig. 4D, right). In the  
 373 next section, we build on these mechanistic characterizations of each regime by examining their  
 374 resilience to optogenetic inactivation.

375 **3.6 EPI inferred SC connectivities reproduce results from optogenetic inacti-  
 376 vation experiments**

377 During the delay period of this task, the circuit must prepare to execute the correct task according  
 378 to the presented cue. The circuit must then maintain a representation of task throughout the delay  
 379 period, which is important for correct execution of the Anti task. Duan et al. found that bilateral  
 380 optogenetic inactivation of SC during the delay period consistently decreased performance in the  
 381 Anti task, but had no effect on the Pro task (Fig. 5A) [48]. The distribution of connectivities  
 382 inferred by EPI exhibited this same effect in simulation at high optogenetic strengths  $\gamma$ , which  
 383 reduce the network activities  $\mathbf{x}(t)$  by a factor  $1 - \gamma$  (Fig. 5B) (see Section 5.5.8).

384 To examine how connectivity affects response to delay period inactivation, we grouped connectivi-  
 385 ties of the EPI distribution along the continuum linking regimes 1 and 2 of Section 3.5:

$$Z(sW) = \{\mathbf{z} \text{ if } \underset{y \in Y}{\operatorname{argmin}} |\mathbf{z} - \mathbf{z}^*(y)|_2 = sW, \text{ for } z \sim q_{\theta}(\mathbf{z} | \mathcal{X})\}, \quad (12)$$

386 where  $Y = \{-1., -0.75, ..., 1.\}$  are the values of  $sW$  for which we calculated the mode.  $Z(sW)$  is  
 387 then the set of EPI samples for which the closest mode was  $\mathbf{z}^*(sW)$  (see Section 5.5.4). In the  
 388 following analyses, we examine how error, and the influence of connectivity eigenvalue on Anti error  
 389 change along this continuum of connectivities. Obtaining the parameter samples for these analysis  
 390 with the learned EPI distribution was more than 20,000 times faster than a brute force approach  
 391 (see Section 5.5.5).

392 The mean increase in Anti error of the EPI distribution is closest to the experimentally measured  
 393 value of 7% at  $\gamma = 0.675$  (Fig. 5B, black dot). At this level of optogenetic strength, regime  
 394 1 exhibits an increase in Anti error with delay period silencing (Fig. 5C, left), while regime 2  
 395 does not. In regime 1, greater  $\lambda_{\text{task}}$  and  $\lambda_{\text{diag}}$  decrease Anti error (Fig. 5C, right). In other words,  
 396 stronger task representations and diagonal amplification make the SC model more resilient to delay  
 397 period silencing in the Anti task. This complements the finding from Duan et al. 2019 [48] that  
 398  $\lambda_{\text{task}}$  and  $\lambda_{\text{diag}}$  improve Anti accuracy.

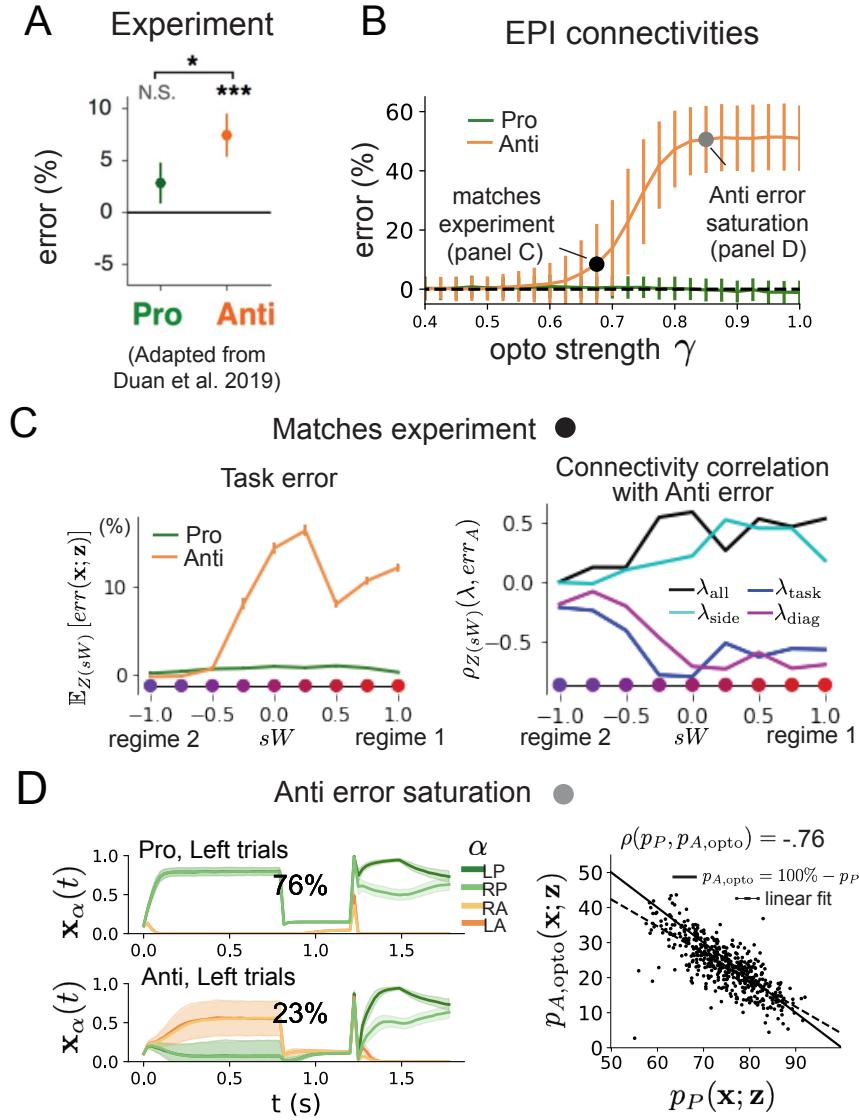


Figure 5: **A.** Mean and standard error (bars) across recording sessions of task error following delay period optogenetic inactivation in rats. **B.** Mean and standard deviation (bars) of task error induced by delay period inactivation of varying optogenetic strength  $\gamma$  across the EPI distribution. **C.** (Left) Mean and standard error of Pro and Anti error from regime 1 to regime 2 at  $\gamma = 0.675$ . (Right) Correlations of connectivity eigenvalues with Anti error from regime 1 to regime 2 at  $\gamma = 0.675$ . **D.** (Left) Mean and standard deviation (shading) of responses of the SC model at the mode of the EPI distribution to delay period inactivation at  $\gamma = 0.85$ . Accuracy in Pro (top) and Anti (bottom) task is shown as a percentage. (Right) Anti accuracy following delay period inactivation at  $\gamma = 0.85$  versus accuracy in the Pro task across connectivities in the EPI distribution.

399 At about  $\gamma = 0.85$  (Fig. 5B, gray dot), the Anti error saturates, while Pro error remains at zero.  
400 Following delay period inactivation at this optogenetic strength, there are strong similarities in  
401 the responses of Pro and Anti trials during the choice period (Fig. 5D, left). We interpreted  
402 these similarities to suggest that delay period inactivation at this saturated level flips the internal  
403 representation of task (from Anti to Pro) in the circuit model. This would explain why the Anti  
404 error saturates at 50%: the average Anti accuracy in EPI inferred connectivities is 75%, but is 25%  
405 when the internal representation is flipped during delay period silencing. This hypothesis prescribes  
406 a model of Anti accuracy during delay period silencing of  $p_{A,\text{opto}} = 100\% - p_P$ , which is fit closely  
407 across both regimes of the EPI inferred connectivities (Fig. 5D, right). Similarities between Pro  
408 and Anti trial responses were not present at the experiment-matching level of  $\gamma = 0.675$  (Fig. S16  
409 left) and neither was anticorrelation in  $p_P$  and  $p_{A,\text{opto}}$  (Fig. S16 right).

410 In summary, the connectivity inferred by EPI to perform rapid task switching replicated results  
411 from optogenetic silencing experiments. We found that at levels of optogenetic strength matching  
412 experimental levels of Anti error, only one regime actually exhibited the effect. This connectivity  
413 regime is less resilient to optogenetic perturbation, and perhaps more biologically realistic. Finally,  
414 we characterized the pathology in Anti error that occurs in both regimes when optogenetic strength  
415 is increased to high levels, leading to a mechanistic hypothesis that is experimentally testable.  
416 The probabilistic tools afforded by EPI yielded this insight: we identified two regimes and the  
417 continuum of connectivities between them by taking gradients of parameter probabilities in the EPI  
418 distribution, we identified sensitivity dimensions by measuring the Hessian of the EPI distribution,  
419 and we obtained many parameter samples at each step along the continuum at an efficient rate.

## 420 4 Discussion

421 In neuroscience, machine learning has primarily been used to reveal structure in neural datasets [20].  
422 Careful inference procedures are developed for these statistical models allowing precise, quantitative  
423 reasoning, which clarifies the way data informs beliefs about the model parameters. However, these  
424 statistical models often lack resemblance to the underlying biology, making it unclear how to go  
425 from the structure revealed by these methods, to the neural mechanisms giving rise to it. In  
426 contrast, theoretical neuroscience has primarily focused on careful models of neural circuits and  
427 the production of emergent properties of computation, rather than measuring structure in neural  
428 datasets. In this work, we improve upon parameter inference techniques in theoretical neuroscience

429 with emergent property inference, harnessing deep learning towards parameter inference in neural  
430 circuit models (see Section 5.1.1).

431 Methodology for statistical inference in circuit models has evolved considerably in recent years.  
432 Early work used rejection sampling techniques [24–26], but EPI and another recently developed  
433 methodology [35] employ deep learning to improve efficiency and provide flexible approximations.  
434 SNPE has been used for posterior inference of parameters in circuit models conditioned upon  
435 exemplar data used to represent computation, but it does not infer parameter distributions that  
436 only produce the computation of interest like EPI (see Section 3.3). When strict control over the  
437 predictions of the inferred parameters is necessary, EPI uses a constrained optimization technique  
438 [38] (see Section 5.1.4) to make inference conditioned on the emergent property possible.

439 A key difference between EPI and SNPE, is that EPI uses gradients of the emergent property  
440 throughout optimization. In Section 3.3, we showed that such gradients confer beneficial scaling  
441 properties, but a concern remains that emergent property gradients may be too computationally  
442 intensive. Even in a case of close biophysical realism with an expensive emergent property gradient,  
443 EPI was run successfully on intermediate hub frequency in a 5-neuron subcircuit model of the STG  
444 (Section 3.1). However, conditioning on the pyloric rhythm [69] in a model of the pyloric subnetwork  
445 model [15] proved to be prohibitive with EPI. The pyloric subnetwork requires many time steps for  
446 simulation and many key emergent property statistics (e.g. burst duration and phase gap) are not  
447 calculable or easily approximated with differentiable functions. In such cases, SNPE, which does  
448 not require differentiability of the emergent property, has proved to be a powerful approach [35].  
449 In summary, choice of deep inference technique should consider emergent property complexity and  
450 differentiability, dimensionality of parameter space, and the importance of constraining the model  
451 behavior predicted by the inferred parameter distribution.

452 In this paper, we prove out the value of deep inference for parameter sensitivity analyses at both the  
453 local and global level. With these techniques, flexible deep probability distributions are optimized  
454 to capture global structure by approximating the full distribution of suitable parameters. Impor-  
455 tantly, the local structure of this deep probability distribution can be quantified at any parameter  
456 choice, offering instant sensitivity measurements after fitting. For example, the global structure  
457 captured by EPI revealed two distinct parameter regimes, which had different local structure quan-  
458 tified by the deep probability distribution (see Section 5.5). In comparison, bayesian MCMC is  
459 considered a popular approach for capturing global parameter structure [68], but there is no vari-  
460 ational approximation (the deep probability distribution in EPI), so sensitivity information is not

461 queryable and sampling remains slow after convergence. Local sensitivity analyses (e.g. [27]) may be  
462 performed independently at individual parameter samples, but these methods alone do not capture  
463 the full picture in nonlinear, complex distributions. In contrast, deep inference yields a probability  
464 distribution that produces a wholistic assessment of parameter sensitivity at the local and global  
465 level, which we used in this study to make novel insights into a range of theoretical models. To-  
466 gether, the abilities to condition upon emergent properties, the efficient inference algorithm, and  
467 the capacity for parameter sensitivity analyses make EPI a useful method for addressing inverse  
468 problems in theoretical neuroscience.

469 **Acknowledgements:**

470 This work was funded by NSF Graduate Research Fellowship, DGE-1644869, McKnight Endow-  
471 ment Fund, NIH NINDS 5R01NS100066, Simons Foundation 542963, NSF NeuroNex Award, DBI-  
472 1707398, The Gatsby Charitable Foundation, Simons Collaboration on the Global Brain Postdoc-  
473 toral Fellowship, Chinese Postdoctoral Science Foundation, and International Exchange Program  
474 Fellowship. Helpful conversations were had with Francesca Mastrogiuseppe, Srdjan Ostojic, James  
475 Fitzgerald, Stephen Baccus, Dhruva Raman, Liam Paninski, and Larry Abbott.

476 **Data availability statement:**

477 The datasets generated during and/or analyzed during the current study are available from the  
478 corresponding author upon reasonable request.

479 **Code availability statement:**

480 All software written for the current study is available at <https://github.com/cunningham-lab/epi>.

481 **References**

- 482 [1] Nancy Kopell and G Bard Ermentrout. Coupled oscillators and the design of central pattern  
483 generators. *Mathematical biosciences*, 90(1-2):87–109, 1988.
- 484 [2] Eve Marder. From biophysics to models of network function. *Annual review of neuroscience*,  
485 21(1):25–45, 1998.
- 486 [3] Larry F Abbott. Theoretical neuroscience rising. *Neuron*, 60(3):489–495, 2008.
- 487 [4] Xiao-Jing Wang. Neurophysiological and computational principles of cortical rhythms in cog-  
488 nition. *Physiological reviews*, 90(3):1195–1268, 2010.

- 489 [5] Timothy O’Leary, Alexander C Sutton, and Eve Marder. Computational models in the age of  
490 large datasets. *Current opinion in neurobiology*, 32:87–94, 2015.
- 491 [6] Ryan N Gutenkunst, Joshua J Waterfall, Fergal P Casey, Kevin S Brown, Christopher R  
492 Myers, and James P Sethna. Universally sloppy parameter sensitivities in systems biology  
493 models. *PLoS Comput Biol*, 3(10):e189, 2007.
- 494 [7] Kamil Erguler and Michael PH Stumpf. Practical limits for reverse engineering of dynamical  
495 systems: a statistical analysis of sensitivity and parameter inferability in systems biology  
496 models. *Molecular BioSystems*, 7(5):1593–1602, 2011.
- 497 [8] Brian K Mannakee, Aaron P Ragsdale, Mark K Transtrum, and Ryan N Gutenkunst. Sloppi-  
498 ness and the geometry of parameter space. In *Uncertainty in Biology*, pages 271–299. Springer,  
499 2016.
- 500 [9] John J Hopfield. Neural networks and physical systems with emergent collective computational  
501 abilities. *Proceedings of the national academy of sciences*, 79(8):2554–2558, 1982.
- 502 [10] Haim Sompolinsky, Andrea Crisanti, and Hans-Jurgen Sommers. Chaos in random neural  
503 networks. *Physical review letters*, 61(3):259, 1988.
- 504 [11] Andrey V Olypher and Ronald L Calabrese. Using constraints on neuronal activity to reveal  
505 compensatory changes in neuronal parameters. *Journal of Neurophysiology*, 98(6):3749–3758,  
506 2007.
- 507 [12] Misha V Tsodyks, William E Skaggs, Terrence J Sejnowski, and Bruce L McNaughton. Para-  
508 doxical effects of external modulation of inhibitory interneurons. *Journal of neuroscience*,  
509 17(11):4382–4388, 1997.
- 510 [13] Kong-Fatt Wong and Xiao-Jing Wang. A recurrent network mechanism of time integration in  
511 perceptual decisions. *Journal of Neuroscience*, 26(4):1314–1328, 2006.
- 512 [14] WR Foster, LH Ungar, and JS Schwaber. Significance of conductances in hodgkin-huxley  
513 models. *Journal of neurophysiology*, 70(6):2502–2518, 1993.
- 514 [15] Astrid A Prinz, Dirk Bucher, and Eve Marder. Similar network activity from disparate circuit  
515 parameters. *Nature neuroscience*, 7(12):1345–1352, 2004.
- 516 [16] Pablo Achard and Erik De Schutter. Complex parameter landscape for a complex neuron  
517 model. *PLoS computational biology*, 2(7):e94, 2006.

- 518 [17] Dmitry Fisher, Itsaso Olasagasti, David W Tank, Emre RF Aksay, and Mark S Goldman.  
519       A modeling framework for deriving the structural and functional architecture of a short-term  
520       memory microcircuit. *Neuron*, 79(5):987–1000, 2013.
- 521 [18] Timothy O’Leary, Alex H Williams, Alessio Franci, and Eve Marder. Cell types, network  
522       homeostasis, and pathological compensation from a biologically plausible ion channel expres-  
523       sion model. *Neuron*, 82(4):809–821, 2014.
- 524 [19] Leandro M Alonso and Eve Marder. Visualization of currents in neural models with similar  
525       behavior and different conductance densities. *Elife*, 8:e42722, 2019.
- 526 [20] Liam Paninski and John P Cunningham. Neural data science: accelerating the experiment-  
527       analysis-theory cycle in large-scale neuroscience. *Current opinion in neurobiology*, 50:232–241,  
528       2018.
- 529 [21] Christopher M Niell and Michael P Stryker. Modulation of visual responses by behavioral state  
530       in mouse visual cortex. *Neuron*, 65(4):472–479, 2010.
- 531 [22] Aman B Saleem, Asli Ayaz, Kathryn J Jeffery, Kenneth D Harris, and Matteo Carandini.  
532       Integration of visual motion and locomotion in mouse visual cortex. *Nature neuroscience*,  
533       16(12):1864–1869, 2013.
- 534 [23] Simon Musall, Matthew T Kaufman, Ashley L Juavinett, Steven Gluf, and Anne K Church-  
535       land. Single-trial neural dynamics are dominated by richly varied movements. *Nature neuro-  
536       science*, 22(10):1677–1686, 2019.
- 537 [24] Mark A Beaumont, Wenyang Zhang, and David J Balding. Approximate bayesian computation  
538       in population genetics. *Genetics*, 162(4):2025–2035, 2002.
- 539 [25] Paul Marjoram, John Molitor, Vincent Plagnol, and Simon Tavaré. Markov chain monte carlo  
540       without likelihoods. *Proceedings of the National Academy of Sciences*, 100(26):15324–15328,  
541       2003.
- 542 [26] Scott A Sisson, Yanan Fan, and Mark M Tanaka. Sequential monte carlo without likelihoods.  
543       *Proceedings of the National Academy of Sciences*, 104(6):1760–1765, 2007.
- 544 [27] Andreas Raue, Clemens Kreutz, Thomas Maiwald, Julie Bachmann, Marcel Schilling, Ursula  
545       Klingmüller, and Jens Timmer. Structural and practical identifiability analysis of partially

- 546 observed dynamical models by exploiting the profile likelihood. *Bioinformatics*, 25(15):1923–  
547 1929, 2009.
- 548 [28] Johan Karlsson, Milena Anguelova, and Mats Jirstrand. An efficient method for structural  
549 identifiability analysis of large dynamic systems. *IFAC Proceedings Volumes*, 45(16):941–946,  
550 2012.
- 551 [29] Keegan E Hines, Thomas R Middendorf, and Richard W Aldrich. Determination of parameter  
552 identifiability in nonlinear biophysical models: A bayesian approach. *Journal of General  
553 Physiology*, 143(3):401–416, 2014.
- 554 [30] Dhruva V Raman, James Anderson, and Antonis Papachristodoulou. Delineating parameter  
555 unidentifiabilities in complex models. *Physical Review E*, 95(3):032314, 2017.
- 556 [31] Gamaleldin F Elsayed and John P Cunningham. Structure in neural population recordings:  
557 an expected byproduct of simpler phenomena? *Nature neuroscience*, 20(9):1310, 2017.
- 558 [32] Cristina Savin and Gašper Tkačik. Maximum entropy models as a tool for building precise  
559 neural controls. *Current opinion in neurobiology*, 46:120–126, 2017.
- 560 [33] Wiktor Mlynarski, Michal Hledík, Thomas R Sokolowski, and Gašper Tkačik. Statistical  
561 analysis and optimality of neural systems. *bioRxiv*, page 848374, 2020.
- 562 [34] Dustin Tran, Rajesh Ranganath, and David Blei. Hierarchical implicit models and likelihood-  
563 free variational inference. In *Advances in Neural Information Processing Systems*, pages 5523–  
564 5533, 2017.
- 565 [35] Pedro J Gonçalves, Jan-Matthis Lueckmann, Michael Deistler, Marcel Nonnenmacher, Kaan  
566 Öcal, Giacomo Bassetto, Chaitanya Chintaluri, William F Podlaski, Sara A Haddad, Tim P  
567 Vogels, et al. Training deep neural density estimators to identify mechanistic models of neural  
568 dynamics. *bioRxiv*, page 838383, 2019.
- 569 [36] Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows.  
570 *International Conference on Machine Learning*, 2015.
- 571 [37] George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji  
572 Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *arXiv preprint  
573 arXiv:1912.02762*, 2019.

- 574 [38] Gabriel Loaiza-Ganem, Yuanjun Gao, and John P Cunningham. Maximum entropy flow  
575 networks. *International Conference on Learning Representations*, 2017.
- 576 [39] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp.  
577 *Proceedings of the 5th International Conference on Learning Representations*, 2017.
- 578 [40] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolu-  
579 tions. In *Advances in neural information processing systems*, pages 10215–10224, 2018.
- 580 [41] Gabrielle J Gutierrez, Timothy O’Leary, and Eve Marder. Multiple mechanisms switch an  
581 electrically coupled, synaptically inhibited neuron between competing rhythmic oscillators.  
582 *Neuron*, 77(5):845–858, 2013.
- 583 [42] Mark S Goldman, Jorge Golowasch, Eve Marder, and LF Abbott. Global structure, robustness,  
584 and modulation of neuronal models. *Journal of Neuroscience*, 21(14):5229–5238, 2001.
- 585 [43] Brendan K Murphy and Kenneth D Miller. Balanced amplification: a new mechanism of  
586 selective amplification of neural activity patterns. *Neuron*, 61(4):635–648, 2009.
- 587 [44] Guillaume Hennequin, Tim P Vogels, and Wulfram Gerstner. Optimal control of transient dy-  
588 namics in balanced networks supports generation of complex movements. *Neuron*, 82(6):1394–  
589 1406, 2014.
- 590 [45] Giulio Bondanelli, Thomas Deneux, Brice Bathellier, and Srdjan Ostojic. Population coding  
591 and network dynamics during off responses in auditory cortex. *BioRxiv*, page 810655, 2019.
- 592 [46] Ashok Litwin-Kumar, Robert Rosenbaum, and Brent Doiron. Inhibitory stabilization and vi-  
593 sual coding in cortical circuits with multiple interneuron subtypes. *Journal of neurophysiology*,  
594 115(3):1399–1409, 2016.
- 595 [47] Agostina Palmigiano, Francesco Fumarola, Daniel P Mossing, Nataliya Kraynyukova, Hillel  
596 Adesnik, and Kenneth Miller. Structure and variability of optogenetic responses identify the  
597 operating regime of cortex. *bioRxiv*, 2020.
- 598 [48] Chunyu A Duan, Marino Pagan, Alex T Piet, Charles D Kopec, Athena Akrami, Alexander J  
599 Riordan, Jeffrey C Erlich, and Carlos D Brody. Collicular circuits for flexible sensorimotor  
600 routing. *bioRxiv*, page 245613, 2019.
- 601 [49] Eve Marder and Vatsala Thirumalai. Cellular, synaptic and network effects of neuromodula-  
602 tion. *Neural Networks*, 15(4-6):479–493, 2002.

- 603 [50] Mark S Goldman. Memory without feedback in a neural network. *Neuron*, 61(4):621–634,  
604 2009.
- 605 [51] Giulio Bondanelli and Srdjan Ostojic. Coding with transient trajectories in recurrent neural  
606 networks. *PLoS computational biology*, 16(2):e1007655, 2020.
- 607 [52] David Sussillo. Neural circuits as computational dynamical systems. *Current opinion in*  
608 *neurobiology*, 25:156–163, 2014.
- 609 [53] Omri Barak. Recurrent neural networks as versatile tools of neuroscience research. *Current*  
610 *opinion in neurobiology*, 46:1–6, 2017.
- 611 [54] Abigail A Russo, Sean R Bittner, Sean M Perkins, Jeffrey S Seely, Brian M London, Antonio H  
612 Lara, Andrew Miri, Najja J Marshall, Adam Kohn, Thomas M Jessell, et al. Motor cortex  
613 embeds muscle-like commands in an untangled population response. *Neuron*, 97(4):953–966,  
614 2018.
- 615 [55] Scott A Sisson, Yanan Fan, and Mark Beaumont. *Handbook of approximate Bayesian compu-*  
616 *tation*. CRC Press, 2018.
- 617 [56] Kyle Cranmer, Johann Brehmer, and Gilles Louppe. The frontier of simulation-based inference.  
618 *Proceedings of the National Academy of Sciences*, 2020.
- 619 [57] Hirofumi Ozeki, Ian M Finn, Evan S Schaffer, Kenneth D Miller, and David Ferster. Inhibitory  
620 stabilization of the cortical network underlies visual surround suppression. *Neuron*, 62(4):578–  
621 592, 2009.
- 622 [58] Daniel B Rubin, Stephen D Van Hooser, and Kenneth D Miller. The stabilized supralinear  
623 network: a unifying circuit motif underlying multi-input integration in sensory cortex. *Neuron*,  
624 85(2):402–417, 2015.
- 625 [59] Guillaume Hennequin, Yashar Ahmadian, Daniel B Rubin, Máté Lengyel, and Kenneth D  
626 Miller. The dynamical regime of sensory cortex: stable dynamics around a single stimulus-  
627 tuned attractor account for patterns of noise variability. *Neuron*, 98(4):846–860, 2018.
- 628 [60] Mark M. Churchland, Byron M. Yu, John P. Cunningham, Leo P. Sugrue, Marlene R. Cohen,  
629 Greg S. Corrado, William T. Newsome, Andrew M. Clark, Paymon Hosseini, Benjamin B.  
630 Scott, David C. Bradley, Matthew A. Smith, Adam Kohn, J. Anthony Movshon, Katherine  
631 M. Armstrong, Tirin Moore, Steve W. Chang, Lawrence H. Snyder, Stephen G. Lisberger,

- 632 Nicholas J. Priebe, Ian M. Finn, David Ferster, Stephen I. Ryu, Gopal Santhanam, Maneesh  
633 Sahani, and Krishna V. Shenoy. Stimulus onset quenches neural variability: a widespread  
634 cortical phenomenon. *Nat. Neurosci.*, 13(3):369–378, 2010.
- 635 [61] Henry Markram, Maria Toledo-Rodriguez, Yun Wang, Anirudh Gupta, Gilad Silberberg, and  
636 Caizhi Wu. Interneurons of the neocortical inhibitory system. *Nature reviews neuroscience*,  
637 5(10):793, 2004.
- 638 [62] Bernardo Rudy, Gordon Fishell, SooHyun Lee, and Jens Hjerling-Leffler. Three groups of  
639 interneurons account for nearly 100% of neocortical gabaergic neurons. *Developmental neuro-*  
640 *biology*, 71(1):45–61, 2011.
- 641 [63] Robin Tremblay, Soohyun Lee, and Bernardo Rudy. GABAergic Interneurons in the Neocortex:  
642 From Cellular Properties to Circuits. *Neuron*, 91(2):260–292, 2016.
- 643 [64] Carsten K Pfeffer, Mingshan Xue, Miao He, Z Josh Huang, and Massimo Scanziani. Inhi-  
644 bition of inhibition in visual cortex: the logic of connections between molecularly distinct  
645 interneurons. *Nature Neuroscience*, 16(8):1068, 2013.
- 646 [65] Daniel J Felleman and David C Van Essen. Distributed hierarchical processing in the primate  
647 cerebral cortex. *Cerebral cortex (New York, NY: 1991)*, 1(1):1–47, 1991.
- 648 [66] C Gardiner. Stochastic methods: A Handbook for the Natural and Social Sciences, 2009.
- 649 [67] Chunyu A Duan, Jeffrey C Erlich, and Carlos D Brody. Requirement of prefrontal and midbrain  
650 regions for rapid executive control of behavior in the rat. *Neuron*, 86(6):1491–1503, 2015.
- 651 [68] Mark Girolami and Ben Calderhead. Riemann manifold langevin and hamiltonian monte  
652 carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*,  
653 73(2):123–214, 2011.
- 654 [69] Eve Marder and Allen I Selverston. *Dynamic biological networks: the stomatogastric nervous*  
655 *system*. MIT press, 1992.
- 656 [70] Lawrence Saul and Michael Jordan. A mean field learning algorithm for unsupervised neural  
657 networks. In *Learning in graphical models*, pages 541–554. Springer, 1998.
- 658 [71] Nicholas Metropolis, Arianna W Rosenbluth, Marshall N Rosenbluth, Augusta H Teller, and  
659 Edward Teller. Equation of state calculations by fast computing machines. *The journal of*  
660 *chemical physics*, 21(6):1087–1092, 1953.

- 661 [72] W Keith Hastings. Monte carlo sampling methods using markov chains and their applications.  
662 1970.
- 663 [73] Ben Calderhead and Mark Girolami. Statistical analysis of nonlinear dynamical systems using  
664 differential geometric sampling methods. *Interface focus*, 1(6):821–835, 2011.
- 665 [74] Andrew Golightly and Darren J Wilkinson. Bayesian parameter inference for stochastic bio-  
666 chemical network models using particle markov chain monte carlo. *Interface focus*, 1(6):807–  
667 820, 2011.
- 668 [75] Oksana A Chkrebtii, David A Campbell, Ben Calderhead, Mark A Girolami, et al. Bayesian  
669 solution uncertainty quantification for differential equations. *Bayesian Analysis*, 11(4):1239–  
670 1267, 2016.
- 671 [76] Juliane Liepe, Paul Kirk, Sarah Filippi, Tina Toni, Chris P Barnes, and Michael PH Stumpf.  
672 A framework for parameter estimation and model selection from experimental data in systems  
673 biology using approximate bayesian computation. *Nature protocols*, 9(2):439–456, 2014.
- 674 [77] Sean R Bittner, Agostina Palmigiano, Kenneth D Miller, and John P Cunningham. Degener-  
675 ate solution networks for theoretical neuroscience. *Computational and Systems Neuroscience  
676 Meeting (COSYNE), Lisbon, Portugal*, 2019.
- 677 [78] Sean R Bittner, Alex T Piet, Chunyu A Duan, Agostina Palmigiano, Kenneth D Miller,  
678 Carlos D Brody, and John P Cunningham. Examining models in theoretical neuroscience with  
679 degenerate solution networks. *Bernstein Conference 2019, Berlin, Germany*, 2019.
- 680 [79] Marcel Nonnenmacher, Pedro J Goncalves, Giacomo Bassetto, Jan-Matthis Lueckmann, and  
681 Jakob H Macke. Robust statistical inference for simulation-based models in neuroscience. In  
682 *Bernstein Conference 2018, Berlin, Germany*, 2018.
- 683 [80] Deistler Michael, , Pedro J Goncalves, Kaan Oecal, and Jakob H Macke. Statistical inference for  
684 analyzing sloppiness in neuroscience models. In *Bernstein Conference 2019, Berlin, Germany*,  
685 2019.
- 686 [81] Jan-Matthis Lueckmann, Pedro J Goncalves, Giacomo Bassetto, Kaan Öcal, Marcel Nonnen-  
687 macher, and Jakob H Macke. Flexible statistical inference for mechanistic models of neural  
688 dynamics. In *Advances in Neural Information Processing Systems*, pages 1289–1299, 2017.

- 689 [82] George Papamakarios, David Sterratt, and Iain Murray. Sequential neural likelihood: Fast  
690 likelihood-free inference with autoregressive flows. In *The 22nd International Conference on*  
691 *Artificial Intelligence and Statistics*, pages 837–848. PMLR, 2019.
- 692 [83] Joeri Hermans, Volodimir Begy, and Gilles Louppe. Likelihood-free mcmc with amortized  
693 approximate ratio estimators. In *International Conference on Machine Learning*, pages 4239–  
694 4248. PMLR, 2020.
- 695 [84] Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and  
696 variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.
- 697 [85] Sean R Bittner and John P Cunningham. Approximating exponential family models (not  
698 single distributions) with a two-network architecture. *arXiv preprint arXiv:1903.07515*, 2019.
- 699 [86] Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary  
700 differential equations. In *Advances in neural information processing systems*, pages 6571–6583,  
701 2018.
- 702 [87] Xuechen Li, Ting-Kam Leonard Wong, Ricky TQ Chen, and David Duvenaud. Scalable  
703 gradients for stochastic differential equations. *arXiv preprint arXiv:2001.01328*, 2020.
- 704 [88] Maria Pia Saccomani, Stefania Audoly, and Leontina D’Angiò. Parameter identifiability of  
705 nonlinear systems: the role of initial conditions. *Automatica*, 39(4):619–632, 2003.
- 706 [89] Stefan Hengl, Clemens Kreutz, Jens Timmer, and Thomas Maiwald. Data-based identifiability  
707 analysis of non-linear dynamical models. *Bioinformatics*, 23(19):2612–2618, 2007.
- 708 [90] George Papamakarios, Theo Pavlakou, and Iain Murray. Masked autoregressive flow for density  
709 estimation. In *Advances in Neural Information Processing Systems*, pages 2338–2347, 2017.
- 710 [91] Durk P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling.  
711 Improved variational inference with inverse autoregressive flow. *Advances in neural information  
712 processing systems*, 29:4743–4751, 2016.
- 713 [92] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *International  
714 Conference on Learning Representations*, 2015.
- 715 [93] Emmanuel Klinger, Dennis Rickert, and Jan Hasenauer. pyabc: distributed, likelihood-free  
716 inference. *Bioinformatics*, 34(20):3591–3593, 2018.

717 [94] David S Greenberg, Marcel Nonnenmacher, and Jakob H Macke. Automatic posterior trans-  
718 formation for likelihood-free inference. *International Conference on Machine Learning*, 2019.

719 [95] Daniel P Mossing, Julia Veit, Agostina Palmigiano, Kenneth D. Miller, and Hillel Adesnik.  
720 Antagonistic inhibitory subnetworks control cooperation and competition across cortical space.  
721 *bioRxiv*, 2021.

722 **5 Methods**

723 **5.1 Emergent property inference (EPI)**

724 Determining the combinations of model parameters that can produce a desired output is a key part  
725 of scientific practice. Solving inverse problems is especially important in neuroscience, since we  
726 require detailed circuit models to produce computation of varying levels of complexity. While much  
727 machine learning research has focused on how to find latent structure in large-scale neural datasets,  
728 less has focused on inverting theoretical circuit models conditioned upon the emergent properties of  
729 computation. Here, we introduce a novel method for statistical inference, which finds distributions  
730 of parameter solutions that are constrained to produce the desired emergent property. This method  
731 seamlessly handles neural circuit models with stochastic nonlinear dynamical generative processes,  
732 which are predominant in theoretical neuroscience.

733 Consider model parameterization  $\mathbf{z}$ , which is a collection of scientifically meaningful variables that  
734 govern the complex simulation of data  $\mathbf{x}$ . For example (see Section 3.1),  $\mathbf{z}$  may be the electrical  
735 conductance parameters of an STG subcircuit, and  $\mathbf{x}$  the evolving membrane potentials (the state)  
736 of the five neurons. In terms of statistical modeling, this circuit model has an intractable likelihood  
737  $p(\mathbf{x} | \mathbf{z})$ , which is predicated by the stochastic differential equations that define the model. From a  
738 theoretical perspective, we are less concerned about the likelihood of an exemplary dataset  $\mathbf{x}$ , but  
739 rather the emergent property of intermediate hub frequency (which implies a consistent dataset  $\mathbf{x}$ ).

740 In the STG example, the statistic  $f(\mathbf{x}; \mathbf{z})$  measures hub neuron frequency from the evolution of  $\mathbf{x}$   
741 governed by parameters  $\mathbf{z}$ . With EPI, we learn distributions of  $\mathbf{z}$  constrained to produce intermediate  
742 hub frequency: to obey the constraints placed on the mean and variance of  $f(\mathbf{x}; \mathbf{z})$ . In general,  
743 an emergent property  $\mathcal{X}$  is defined through the choice of  $f(\mathbf{x}; \mathbf{z})$  (which may be one or multiple  
744 statistics), and its means  $\boldsymbol{\mu}$ , and variances  $\boldsymbol{\sigma}^2$ :

$$\mathcal{X} : \mathbb{E}_{\mathbf{z}, \mathbf{x}} [f(\mathbf{x}; \mathbf{z})] = \boldsymbol{\mu}, \text{Var}_{\mathbf{z}, \mathbf{x}} [f(\mathbf{x}; \mathbf{z})] = \boldsymbol{\sigma}^2. \quad (13)$$

745 Precisely, the emergent property statistics  $f(\mathbf{x}; \mathbf{z})$  must have means  $\boldsymbol{\mu}$  and variances  $\boldsymbol{\sigma}^2$  over the EPI  
746 distribution of parameters and the data produced by those parameters. Technically, an emergent  
747 property may be a combination of first-, second-, or higher-order moments, but this study focuses  
748 on the case written in Equation 13.

749 In EPI, deep probability distributions are optimized to learn the inferred distribution. In deep  
750 probability distributions, a simple random variable  $\mathbf{z}_0 \sim q_0(\mathbf{z}_0)$  (we choose an isotropic gaussian)

751 is mapped deterministically via a sequence of deep neural network layers ( $g_1, \dots g_l$ ) parameterized  
 752 by weights and biases  $\theta$  to the support of the distribution of interest:

$$\mathbf{z} = g_\theta(\mathbf{z}_0) = g_l(\dots g_1(\mathbf{z}_0)) \sim q_\theta(\mathbf{z}). \quad (14)$$

753 Such deep probability distributions embed the inferred distribution in a deep network. Once op-  
 754 timized, this deep network representation has remarkably useful properties: fast sampling and  
 755 probability evaluations Importantly, fast probability evaluations confer fast gradient and Hessian  
 756 calculations as well.

757 Given this choice of circuit model and emergent property  $\mathcal{X}$ ,  $q_\theta(\mathbf{z})$  is optimized via the neural  
 758 network parameters  $\theta$  to find a maximally entropic distribution  $q_\theta^*$  within the deep variational  
 759 family  $\mathcal{Q}$  producing the emergent property  $\mathcal{X}$ :

$$q_\theta(\mathbf{z} | \mathcal{X}) = q_\theta^*(\mathbf{z}) = \operatorname{argmax}_{q_\theta \in \mathcal{Q}} H(q_\theta(\mathbf{z})) \quad (15)$$

s.t.  $\mathcal{X} : \mathbb{E}_{\mathbf{z}, \mathbf{x}} [f(\mathbf{x}; \mathbf{z})] = \boldsymbol{\mu}, \operatorname{Var}_{\mathbf{z}, \mathbf{x}} [f(\mathbf{x}; \mathbf{z})] = \boldsymbol{\sigma}^2$ .

760 Entropy is chosen as the normative selection principle to match that of variational bayesian methods  
 761 (see Section 5.1.3). However, a key difference is that variational bayesian methods do not constrain  
 762 the predictions of their inferred parameter distribution. This optimization is executed using the  
 763 algorithm of Maximum Entropy Flow Networks (MEFNs) [38].

764 In the remainder of Section 5.1, we will explain the finer details and motivation of the EPI method.  
 765 First, we explain related approaches and what EPI introduces to this domain (Section 5.1.1). Sec-  
 766 ond, we describe the special class of deep probability distributions used in EPI called normalizing  
 767 flows (Section 5.1.2). Then, we establish the known relationship between maximum entropy dis-  
 768 tributions and exponential families (Section 5.1.3). Next, we explain the constrained optimization  
 769 technique used to solve Equation 15 (Section 5.1.4). Then, we demonstrate the details of this  
 770 optimization in a toy example (Section 5.1.5). Finally, we explain how EPI is a form of variational  
 771 inference (Section 5.1.6).

### 772 5.1.1 Related approaches

773 When bayesian inference problems lack conjugacy, scientists use approximate inference methods like  
 774 variational inference (VI) [70] and Markov chain Monte Carlo (MCMC) [71, 72]. After optimization,  
 775 variational methods return a parameterized posterior distribution, which we can analyze. Also, the  
 776 variational approximating distribution class is often chosen such that it permits fast sampling. In

777 contrast MCMC methods only produce samples from the approximated posterior distribution. No  
778 parameterized distribution is estimated, and additional samples are always generated with the same  
779 sampling complexity. Inference in models defined by systems of differential has been demonstrated  
780 with MCMC [68], although this approach requires tractable likelihoods. Advancements have intro-  
781 duced sampling [73], likelihood approximation [74], and uncertainty quantification techniques [75]  
782 to make MCMC approaches more efficient and expand the class of applicable models.

783 Simulation-based inference [56] is model parameter inference in the absence of a tractable likeli-  
784 hood function. The most prevalent approach to simulation-based inference is approximate bayesian  
785 computation (ABC) [24], in which satisfactory parameter samples are kept from random prior sam-  
786 pling according to a rejection heuristic. The obtained set of parameters do not have a probabilities,  
787 and further insight about the model must be gained from examination of the parameter set and  
788 their generated activity. Methodological advances to ABC methods have come through the use of  
789 Markov chain Monte Carlo (MCMC-ABC) [25] and sequential Monte Carlo (SMC-ABC) [26] sam-  
790 pling techniques. SMC-ABC is considered state-of-the-art ABC, yet this approach still struggles  
791 to scale in dimensionality [55] (cf. Fig. 2). Still, this method has enjoyed much success in systems  
792 biology [76]. Furthermore, once a parameter set has been obtained by SMC-ABC from a finite set  
793 of particles, the SMC-ABC algorithm must be run again from scratch with a new population of  
794 initialized particles to obtain additional samples.

795 For scientific model analysis, we seek a parameter distribution represented by an approximating  
796 distribution as in variational inference [70]: a variational approximation that once optimized yields  
797 fast analytic calculations and samples. For the reasons described above, ABC and MCMC tech-  
798 niques are unattractive, since they only produce a set of parameter samples lacking probabilities  
799 and have unchanging sampling rate. EPI infers parameters in circuit models using the MEFN [38]  
800 algorithm with a deep variational approximation. The deep neural network of EPI (Fig. 1E) de-  
801 fines the parametric form (with variational parameters  $\theta$ ) of the deep variational approximation of  
802 circuit parameters  $\mathbf{z}$ .

803 Since EPI is not conditioning upon exemplary data as in variational bayesian methodology, EPI  
804 is not doing established variational inference. In contrast, the EPI distribution is constrained to  
805 produce an emergent property. EPI optimization is enabled using stochastic gradient techniques in  
806 the spirit of likelihood-free variational inference [34]. The analytic relationship between EPI and  
807 variational inference is explained in Section 5.1.6.

808 We note that, during our preparation and early presentation of this work [77, 78], another work

809 has arisen with broadly similar goals: bringing statistical inference to mechanistic models of neural  
810 circuits ([35, 79, 80]). We are encouraged by this general problem being recognized by others in the  
811 community, and we emphasize that these works offer complementary neuroscientific contributions  
812 (different theoretical models of focus) and use different technical methodologies (ours is built on  
813 our prior work [38], theirs similarly [81]).

814 The method EPI differs from SNPE in some key ways. SNPE belongs to a “sequential” class  
815 of recently developed simulation-based inference methods in which two neural networks are used  
816 for posterior inference. This first neural network is a deep probability distribution (normalizing  
817 flow) used to estimate the posterior  $p(\mathbf{z} | \mathbf{x})$  (SNPE) or the likelihood  $p(\mathbf{x} | \mathbf{z})$  (sequential neural  
818 likelihood (SNL [82])). A recent advance uses an unconstrained neural network to estimate the  
819 likelihood ratio (sequential neural ratio estimation (SNRE [83])). In SNL and SNRE, MCMC  
820 sampling techniques are used to obtain samples from the approximated posterior. This contrasts  
821 with EPI and SNPE, which use deep probability distributions to model parameters, which facilitates  
822 immediate measurements of sample probability, gradient, or Hessian for system analysis. The  
823 second neural network in this sequential class of methods is the amortizer. This unconstrained  
824 deep network maps data  $\mathbf{x}$  (or statistics  $f(\mathbf{x}; \mathbf{z})$ ) or model parameters  $\mathbf{z}$  to the weights and biases of  
825 the first neural network. These methods are optimized on a conditional density (or ratio) estimation  
826 objective. The data used to optimize this objective are generated via an adaptive procedure, in  
827 which training data pairs  $(\mathbf{x}_i, \mathbf{z}_i)$  become sequentially closer to the true data and posterior.

828 The approximating fidelity of the deep probability distribution in sequential approaches is opti-  
829 mized to generalize across the training distribution of the conditioning variable. This generalization  
830 property of the sequential methods can reduce the accuracy at the singular posterior of interest.  
831 Whereas in EPI, the entire expressivity of the deep probability distribution is dedicated to learning  
832 a single distribution as well as possible. Amortization is not possible in EPI, since EPI learns  
833 an exponential family distribution parameterized by its mean (see Section 5.1.3). Since EPI dis-  
834 tributions are defined by the mean  $\mu$  of their statistics, there is the well-known inverse mapping  
835 problem of exponential families [84] that prohibits an amortization based approach. However, we  
836 have shown that the same two-network architecture of the sequential simulation-based inference  
837 methods can be used for amortized inference in intractable exponential family posteriors using their  
838 natural parameterization [85].

839 Finally, one important differentiating factor between EPI and sequential simulation-based infer-  
840 ence methods is that EPI leverages gradients  $\nabla_{\mathbf{z}} f(\mathbf{x}; \mathbf{z})$  during optimization. These gradients can

841 improve convergence time and scalability, as we have shown on an example conditioning low-rank  
 842 RNN connectivity on the property of stable amplification (see Section 3.3). With EPI, we prove  
 843 out the suggestion that a deep inference technique can improve efficiency by leveraging these model  
 844 gradients when they are tractable. Sequential simulation-based inference techniques may be better  
 845 suited for scientific problems where  $\nabla_{\mathbf{z}} f(\mathbf{x}; \mathbf{z})$  is intractable or unavailable, like when there is a non-  
 846 differentiable model or it requires lengthy simulations. However, the sequential simulation-based  
 847 inference techniques cannot constrain the predictions of the inferred distribution in the manner of  
 848 EPI.

849 Structural identifiability analysis involves the measurement of sensitivity and unidentifiabilities in  
 850 scientific models. Around a single parameter choice, one can measure the Jacobian. One approach  
 851 for this calculation that scales well is EAR [28]. A popular efficient approach for systems of ODEs  
 852 has been neural ODE adjoint [86] and its stochastic adaptation [87]. Casting identifiability as a  
 853 statistical estimation problem, the profile likelihood works via iterated optimization while holding  
 854 parameters fixed [27]. An exciting recent method is capable of recovering the functional form of such  
 855 unidentifiabilities away from a point by following degenerate dimensions of the fisher information  
 856 matrix [30]. Global structural non-identifiabilities can be found for models with polynomial or  
 857 rational dynamics equations using DAISY [88], or through mean optimal transformations [89].  
 858 With EPI, we have all the benefits given by a statistical inference method plus the ability to query  
 859 the first- or second-order gradient of the probability of the inferred distribution at any chosen  
 860 parameter value. The second-order gradient of the log probability (the Hessian), which is directly  
 861 afforded by EPI distributions, produces quantified information about parametric sensitivity of the  
 862 emergent property in parameter space (see Section 3.2).

### 863 5.1.2 Deep probability distributions and normalizing flows

864 Deep probability distributions are comprised of multiple layers of fully connected neural networks  
 865 (Equation 14). When each neural network layer is restricted to be a bijective function, the sample  
 866 density can be calculated using the change of variables formula at each layer of the network. For  
 867  $\mathbf{z}_i = g_i(\mathbf{z}_{i-1})$ ,

$$p(\mathbf{z}_i) = p(g_i^{-1}(\mathbf{z}_i)) \left| \det \frac{\partial g_i^{-1}(\mathbf{z}_i)}{\partial \mathbf{z}_i} \right| = p(\mathbf{z}_{i-1}) \left| \det \frac{\partial g_i(\mathbf{z}_{i-1})}{\partial \mathbf{z}_{i-1}} \right|^{-1}. \quad (16)$$

868 However, this computation has cubic complexity in dimensionality for fully connected layers. By  
 869 restricting our layers to normalizing flows [36, 37] – bijective functions with fast log determinant

870 Jacobian computations, which confer a fast calculation of the sample log probability. Fast log  
871 probability calculation confers efficient optimization of the maximum entropy objective (see Section  
872 5.1.4).

873 We use the Real NVP [39] normalizing flow class, because its coupling architecture confers both  
874 fast sampling (forward) and fast log probability evaluation (backward). Fast probability evaluation  
875 facilitates fast gradient and Hessian evaluation of log probability throughout parameter space.  
876 Glow permutations were used in between coupling stages [40]. This is in contrast to autoregressive  
877 architectures [90, 91], in which only one of the forward or backward passes can be efficient. In this  
878 work, normalizing flows are used as flexible parameter distribution approximations  $q_{\theta}(\mathbf{z})$  having  
879 weights and biases  $\theta$ . We specify the architecture used in each application by the number of Real-  
880 NVP affine coupling stages, and the number of neural network layers and units per layer of the  
881 conditioning functions.

882 When calculating Hessians of log probabilities in deep probability distributions, it is important to  
883 consider the normalizing flow architecture. With autoregressive architectures [90, 91], fast sam-  
884 pling and fast log probability evaluations are mutually exclusive. That makes these architectures  
885 undesirable for EPI, where efficient sampling is important for optimization, and log probability  
886 evaluation speed predicates the efficiency of gradient and Hessian calculations. With Real NVP  
887 coupling architectures, we get both fast sampling and fast Hessians making both optimization and  
888 scientific analysis efficient.

### 889 5.1.3 Maximum entropy distributions and exponential families

890 EPI is a maximum entropy distribution, which have fundamental links to exponential family dis-  
891 tributions. A maximum entropy distribution of form:

$$p^*(\mathbf{z}) = \underset{p \in \mathcal{P}}{\operatorname{argmax}} H(p(\mathbf{z})) \quad (17)$$

s.t.  $\mathbb{E}_{\mathbf{z} \sim p}[T(\mathbf{z})] = \boldsymbol{\mu}_{\text{opt}}$ .

892 will have probability density in the exponential family:

$$p^*(\mathbf{z}) \propto \exp(\boldsymbol{\eta}^\top T(\mathbf{z})). \quad (18)$$

893 The mappings between the mean parameterization  $\boldsymbol{\mu}_{\text{opt}}$  and the natural parameterization  $\boldsymbol{\eta}$  are  
894 formally hard to identify except in special cases [84].

895 In EPI, emergent properties are defined as statistics having a fixed mean and variance as in Equations 2 and 3. The variance constraint is a second moment constraint on  $f(\mathbf{x}; \mathbf{z})$

$$\text{Var}_{\mathbf{z}, \mathbf{x}} [f(\mathbf{x}; \mathbf{z})] = \mathbb{E}_{\mathbf{z}, \mathbf{x}} [(f(\mathbf{x}; \mathbf{z}) - \boldsymbol{\mu})^2] \quad (19)$$

897 As a general maximum entropy distribution (Equation 17), the sufficient statistics vector contains  
898 both first and second order moments of  $f(\mathbf{x}; \mathbf{z})$

$$T(\mathbf{x}; \mathbf{z}) = \begin{bmatrix} f(\mathbf{x}; \mathbf{z}) \\ (f(\mathbf{x}; \mathbf{z}) - \boldsymbol{\mu})^2 \end{bmatrix}, \quad (20)$$

899 which are constrained to the chosen means and variances

$$\boldsymbol{\mu}_{\text{opt}} = \begin{bmatrix} \boldsymbol{\mu} \\ \sigma^2 \end{bmatrix}. \quad (21)$$

#### 900 5.1.4 Augmented lagrangian optimization

901 To optimize  $q_{\boldsymbol{\theta}}(\mathbf{z})$  in Equation 15, the constrained maximum entropy optimization is executed using  
902 the augmented lagrangian method. The following objective is minimized:

$$L(\boldsymbol{\theta}; \boldsymbol{\eta}_{\text{opt}}, c) = -H(q_{\boldsymbol{\theta}}) + \boldsymbol{\eta}_{\text{opt}}^\top R(\boldsymbol{\theta}) + \frac{c}{2} \|R(\boldsymbol{\theta})\|^2 \quad (22)$$

903 where average constraint violations  $R(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\theta}}} [\mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}|\mathbf{z})} [T(\mathbf{x}; \mathbf{z}) - \boldsymbol{\mu}_{\text{opt}}]]$ ,  $\boldsymbol{\eta}_{\text{opt}} \in \mathbb{R}^m$  are the  
904 Lagrange multipliers where  $m = |\boldsymbol{\mu}_{\text{opt}}| = |T(\mathbf{x}; \mathbf{z})| = 2|f(\mathbf{x}; \mathbf{z})|$ , and  $c$  is the penalty coefficient. The  
905 sufficient statistics  $T(\mathbf{x}; \mathbf{z})$  and mean parameter  $\boldsymbol{\mu}_{\text{opt}}$  are determined by the means  $\boldsymbol{\mu}$  and variances  
906  $\sigma^2$  of emergent property statistics  $f(\mathbf{x}; \mathbf{z})$  defined in Equation 15 (see Section 5.1.6). Specifically,  
907  $T(\mathbf{x}; \mathbf{z})$  is a concatenation of the first and second moments,  $\boldsymbol{\mu}_{\text{opt}}$  is a concatenation of  $\boldsymbol{\mu}$  and  $\sigma^2$   
908 (see section 5.1.3), and the Lagrange multipliers are closely related to the natural parameters  $\boldsymbol{\eta}$  of  
909 exponential families (see Section 5.1.3). Weights and biases  $\boldsymbol{\theta}$  of the deep probability distribution  
910 are optimized according to Equation 22 using the Adam optimizer with learning rate  $10^{-3}$  [92].

911 The gradient with respect to entropy  $H(q_{\boldsymbol{\theta}}(\mathbf{z}))$  can be expressed using the reparameterization trick  
912 as an expectation of the negative log density of parameter samples  $\mathbf{z}$  over the randomness in the  
913 parameterless initial distribution  $q_0(\mathbf{z}_0)$ :

$$H(q_{\boldsymbol{\theta}}(\mathbf{z})) = \int -q_{\boldsymbol{\theta}}(\mathbf{z}) \log(q_{\boldsymbol{\theta}}(\mathbf{z})) d\mathbf{z} = \mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\theta}}} [-\log(q_{\boldsymbol{\theta}}(\mathbf{z}))] = \mathbb{E}_{\mathbf{z}_0 \sim q_0} [-\log(q_{\boldsymbol{\theta}}(g_{\boldsymbol{\theta}}(\mathbf{z}_0)))]. \quad (23)$$

914 Thus, the gradient of the entropy of the deep probability distribution can be estimated as an  
915 average with respect to the base distribution  $\mathbf{z}_0$ :

$$\nabla_{\boldsymbol{\theta}} H(q_{\boldsymbol{\theta}}(\mathbf{z})) = \mathbb{E}_{\mathbf{z}_0 \sim q_0} [-\nabla_{\boldsymbol{\theta}} \log(q_{\boldsymbol{\theta}}(g_{\boldsymbol{\theta}}(\mathbf{z}_0)))]. \quad (24)$$

916 The lagrangian parameters  $\eta_{\text{opt}}$  are initialized to zero and adapted following each augmented  
917 lagrangian epoch, which is a period of optimization with fixed  $(\eta_{\text{opt}}, c)$  for a given number of  
918 stochastic optimization iterations. A low value of  $c$  is used initially, and conditionally increased  
919 after each epoch based on constraint error reduction. The penalty coefficient is updated based  
920 on the result of a hypothesis test regarding the reduction in constraint violation. The p-value of  
921  $\mathbb{E}[|R(\theta_{k+1})|] > \gamma \mathbb{E}[|R(\theta_k)|]$  is computed, and  $c_{k+1}$  is updated to  $\beta c_k$  with probability  $1 - p$ . The  
922 other update rule is  $\eta_{\text{opt},k+1} = \eta_{\text{opt},k} + c_k \frac{1}{n} \sum_{i=1}^n (T(\mathbf{x}^{(i)}) - \mu_{\text{opt}})$  given a batch size  $n$ . Throughout  
923 the study,  $\gamma = 0.25$ , while  $\beta$  was chosen to be either 2 or 4. The batch size of EPI also varied  
924 according to application.

925 The intention is that  $c$  and  $\eta_{\text{opt}}$  start at values encouraging entropic growth early in optimization.  
926 With each training epoch in which the update rule for  $c$  is invoked by unsatisfactory constraint  
927 error reduction, the constraint satisfaction terms are increasingly weighted, resulting in a decreased  
928 entropy. This encourages the discovery of suitable regions of parameter space, and the subsequent  
929 refinement of the distribution to produce the emergent property (see example in Section 5.1.5). The  
930 momentum parameters of the Adam optimizer are reset at the end of each augmented lagrangian  
931 epoch.

932 Rather than starting optimization from some  $\theta$  drawn from a randomized distribution, we found  
933 that initializing  $q_{\theta}(\mathbf{z})$  to approximate an isotropic Gaussian distribution conferred more stable, con-  
934 sistent optimization. The parameters of the Gaussian initialization were chosen on an application-  
935 specific basis. Throughout the study, we chose isotropic Gaussian initializations with mean  $\mu_{\text{init}}$   
936 at the center of the distribution support and some standard deviation  $\sigma_{\text{init}}$ , except for one case,  
937 where an initialization informed by random search was used (see Section 5.2).

938 To assess whether the EPI distribution  $q_{\theta}(\mathbf{z})$  produces the emergent property, we assess whether  
939 each individual constraint on the means and variances of  $f(\mathbf{x}; \mathbf{z})$  is satisfied. We consider the EPI  
940 to have converged when a null hypothesis test of constraint violations  $R(\theta)_i$  being zero is accepted  
941 for all constraints  $i \in \{1, \dots, m\}$  at a significance threshold  $\alpha = 0.05$ . This significance threshold is  
942 adjusted through Bonferroni correction according to the number of constraints  $m$ . The p-values for  
943 each constraint are calculated according to a two-tailed nonparametric test, where 200 estimations  
944 of the sample mean  $R(\theta)^i$  are made using  $N_{\text{test}}$  samples of  $\mathbf{z} \sim q_{\theta}(\mathbf{z})$  at the end of the augmented  
945 lagrangian epoch.

946 When assessing the suitability of EPI for a particular modeling question, there are some important  
947 technical considerations. First and foremost, as in any optimization problem, the defined emergent

948 property should always be appropriately conditioned (constraints should not have wildly different  
 949 units). Furthermore, if the program is underconstrained (not enough constraints), the distribution  
 950 grows (in entropy) unstably unless mapped to a finite support. If overconstrained, there is no pa-  
 951 rameter set producing the emergent property, and EPI optimization will fail (appropriately). Next,  
 952 one should consider the computational cost of the gradient calculations. In the best circumstance,  
 953 there is a simple, closed form expression (e.g. Section 5.3) for the emergent property statistic given  
 954 the model parameters. On the other end of the spectrum, many forward simulation iterations  
 955 may be required before a high quality measurement of the emergent property statistic is available  
 956 (e.g. Section 5.2). In such cases, backpropagating gradients through the SDE evolution will be  
 957 expensive.

### 958 5.1.5 Example: 2D LDS

959 To gain intuition for EPI, consider a two-dimensional linear dynamical system (2D LDS) model  
 960 (Fig. S1A):

$$961 \quad \tau \frac{d\mathbf{x}}{dt} = A\mathbf{x} \quad (25)$$

961 with

$$962 \quad A = \begin{bmatrix} a_1 & a_2 \\ a_3 & a_4 \end{bmatrix}. \quad (26)$$

962 To run EPI with the dynamics matrix elements as the free parameters  $\mathbf{z} = [a_1, a_2, a_3, a_4]$  (fixing  
 963  $\tau = 1s$ ), the emergent property statistics  $f(\mathbf{x}; \mathbf{z})$  were chosen to contain the oscillatory frequency,  
 964  $\frac{\text{imag}(\lambda_1)}{2\pi}$ , and the growth/decay factor,  $\text{real}(\lambda_1)$ , of the oscillating system.  $\lambda_1$  is the eigenvalue of  
 965 greatest real part when the imaginary component is zero, and alternatively of positive imaginary  
 966 component when the eigenvalues are complex conjugate pairs. To learn the distribution of real  
 967 entries of  $A$  that produce a band of oscillating systems around 1Hz, we formalized this emergent  
 968 property as  $\text{real}(\lambda_1)$  having mean zero with variance  $0.25^2$ , and the oscillation frequency  $2\pi\text{imag}(\lambda_1)$   
 969 having mean 1Hz with variance  $(0.1\text{Hz})^2$ :

$$970 \quad \mathbb{E}[T(\mathbf{x})]_{\mathbf{z}, \mathbf{x}} \triangleq \mathbb{E} \begin{bmatrix} \text{real}(\lambda_1)(\mathbf{x}; \mathbf{z}) \\ \text{imag}(\lambda_1)(\mathbf{x}; \mathbf{z}) \\ (\text{real}(\lambda_1)(\mathbf{x}; \mathbf{z}) - 0)^2 \\ (\text{imag}(\lambda_1)(\mathbf{x}; \mathbf{z}) - 2\pi\omega)^2 \end{bmatrix} = \begin{bmatrix} 0.0 \\ 2\pi \\ 0.25^2 \\ (2\pi 0.1)^2 \end{bmatrix} \triangleq \boldsymbol{\mu}_{\text{opt.}} \quad (27)$$

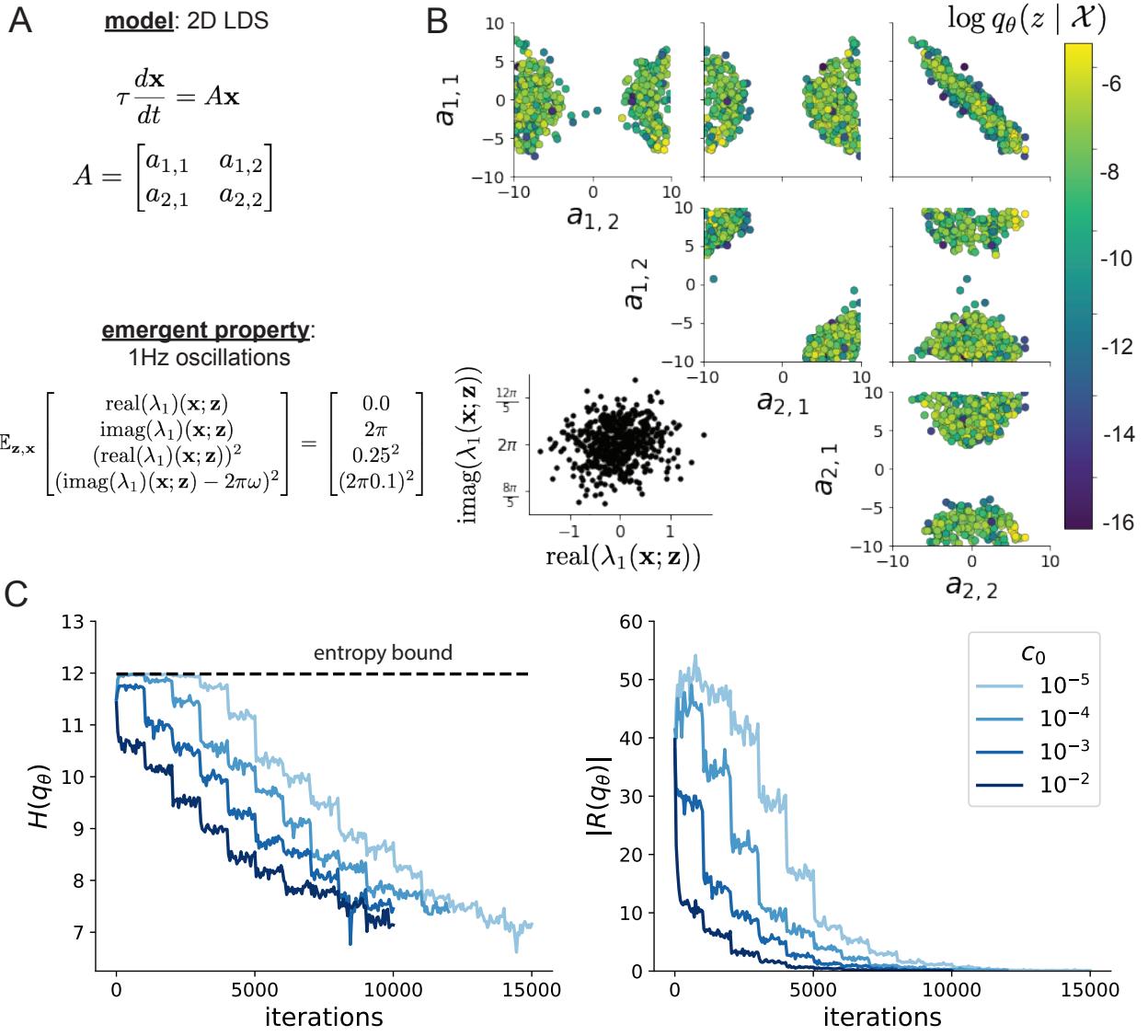


Figure S1: **A.** Two-dimensional linear dynamical system model, where real entries of the dynamics matrix  $A$  are the parameters. **B.** The EPI distribution for a two-dimensional linear dynamical system with  $\tau = 1$  that produces an average of 1Hz oscillations with some small amount of variance. Dashed lines indicate the parameter axes. **C.** Entropy throughout the optimization. At the beginning of each augmented lagrangian epoch (2,000 iterations), the entropy dipped due to the shifted optimization manifold where emergent property constraint satisfaction is increasingly weighted. **D.** Emergent property moments throughout optimization. At the beginning of each augmented lagrangian epoch, the emergent property moments adjust closer to their constraints.

971 Unlike the models we presented in the main text, this model admits an analytical form for the  
 972 mean emergent property statistics given parameter  $\mathbf{z}$ , since the eigenvalues can be calculated using  
 973 the quadratic formula:

$$\lambda = \frac{\left(\frac{a_1+a_4}{\tau}\right) \pm \sqrt{\left(\frac{a_1+a_4}{\tau}\right)^2 + 4\left(\frac{a_2a_3-a_1a_4}{\tau}\right)}}{2}. \quad (28)$$

974 Importantly, even though  $\mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}|\mathbf{z})}[T(\mathbf{x})]$  is calculable directly via a closed form function and  
 975 does not require simulation, we cannot derive the distribution  $q_{\theta}^*$  directly. This fact is due to the  
 976 formally hard problem of the backward mapping: finding the natural parameters  $\eta$  from the mean  
 977 parameters  $\mu$  of an exponential family distribution [84]. Instead, we used EPI to approximate this  
 978 distribution (Fig. S1B). We used a real-NVP normalizing flow architecture with four masks, two  
 979 neural network layers of 15 units per mask, with batch normalization momentum 0.99, mapped  
 980 onto a support of  $z_i \in [-10, 10]$ . (see Section 5.1.2).

981 Even this relatively simple system has nontrivial (though intuitively sensible) structure in the  
 982 parameter distribution. To validate our method, we analytically derived the contours of the prob-  
 983 ability density from the emergent property statistics and values. In the  $a_1$ - $a_4$  plane, the black  
 984 line at  $\text{real}(\lambda_1) = \frac{a_1+a_4}{2} = 0$ , dotted black line at the standard deviation  $\text{real}(\lambda_1) = \frac{a_1+a_4}{2} \pm 0.25$ ,  
 985 and the dotted gray line at twice the standard deviation  $\text{real}(\lambda_1) = \frac{a_1+a_4}{2} \pm 0.5$  follow the contour  
 986 of probability density of the samples (Fig. S2A). The distribution precisely reflects the desired  
 987 statistical constraints and model degeneracy in the sum of  $a_1$  and  $a_4$ . Intuitively, the parameters  
 988 equivalent with respect to emergent property statistic  $\text{real}(\lambda_1)$  have similar log densities.

989 To explain the bimodality of the EPI distribution, we examined the imaginary component of  $\lambda_1$ .  
 990 When  $\text{real}(\lambda_1) = \frac{a_1+a_4}{2} = 0$ , we have

$$\text{imag}(\lambda_1) = \begin{cases} \sqrt{\frac{a_1a_4-a_2a_3}{\tau}}, & \text{if } a_1a_4 < a_2a_3 \\ 0 & \text{otherwise} \end{cases}. \quad (29)$$

991 When  $\tau = 1$  and  $a_1a_4 > a_2a_3$  (center of distribution above), we have the following equation for the  
 992 other two dimensions:

$$\text{imag}(\lambda_1)^2 = a_1a_4 - a_2a_3 \quad (30)$$

993 Since we constrained  $\mathbb{E}_{\mathbf{z} \sim q_{\theta}}[\text{imag}(\lambda)] = 2\pi$ , we can plot contours of the equation  $\text{imag}(\lambda_1)^2 =$   
 994  $a_1a_4 - a_2a_3 = (2\pi)^2$  for various  $a_1a_4$  (Fig. S2B). With  $\sigma_{1,4} = \mathbb{E}_{\mathbf{z} \sim q_{\theta}}(|a_1a_4 - E_{q_{\theta}}[a_1a_4]|)$ , we show  
 995 the contours as  $a_1a_4 = 0$  (black),  $a_1a_4 = -\sigma_{1,4}$  (black dotted), and  $a_1a_4 = -2\sigma_{1,4}$  (grey dotted).  
 996 This validates the curved structure of the inferred distribution learned through EPI. We took steps

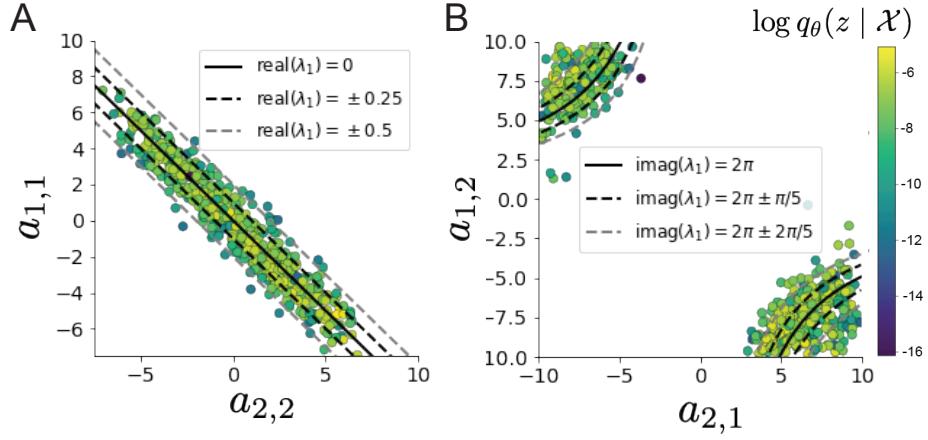


Figure S2: **A.** Probability contours in the  $a_1-a_4$  plane were derived from the relationship to emergent property statistic of growth/decay factor  $\text{real}(\lambda_1)$ . **B.** Probability contours in the  $a_2-a_3$  plane were derived from the emergent property statistic of oscillation frequency  $2\pi\text{imag}(\lambda_1)$ .

in negative standard deviation of  $a_1a_4$  (dotted and gray lines), since there are few positive values  $a_1a_4$  in the learned distribution. Subtler combinations of model and emergent property will have more complexity, further motivating the use of EPI for understanding these systems. As we expect, the distribution results in samples of two-dimensional linear systems oscillating near 1Hz (Fig. S3).

### 5.1.6 EPI as variational inference

In bayesian inference a prior belief about model parameters  $\mathbf{z}$  is stated in a prior distribution  $p(\mathbf{z})$ , and the statistical model capturing the effect of  $\mathbf{z}$  on observed data points  $\mathbf{x}$  is formalized in the likelihood distribution  $p(\mathbf{x} | \mathbf{z})$ . In bayesian inference, we obtain a posterior distribution  $p(\mathbf{z} | \mathbf{x})$ , which captures how the data inform our knowledge of model parameters using Bayes' rule:

$$p(\mathbf{z} | \mathbf{x}) = \frac{p(\mathbf{x} | \mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}. \quad (31)$$

The posterior distribution is analytically available when the prior is conjugate with the likelihood. However, conjugacy is rare in practice, and alternative methods, such as variational inference [70], are utilized.

In variational inference, a posterior approximation  $q_{\theta}^*$  is chosen from within some variational family  $\mathcal{Q}$

$$q_{\theta}^*(\mathbf{z}) = \underset{q_{\theta} \in \mathcal{Q}}{\operatorname{argmin}} KL(q_{\theta}(\mathbf{z}) || p(\mathbf{z} | \mathbf{x})). \quad (32)$$

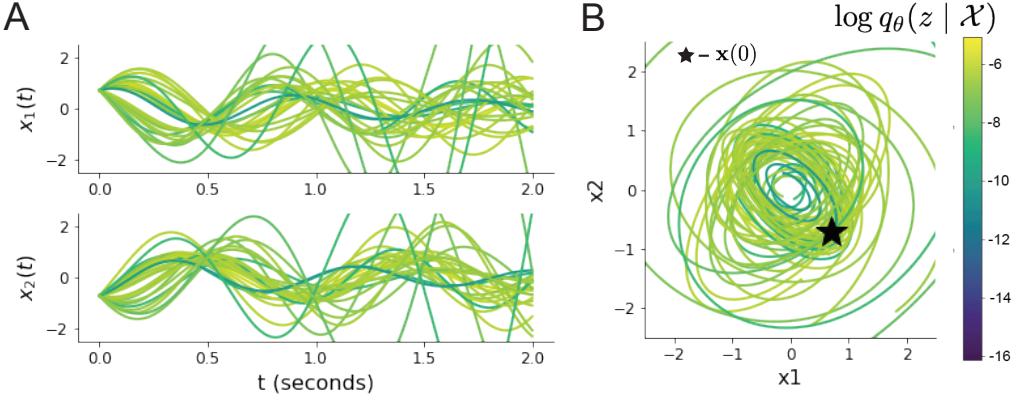


Figure S3: Sampled dynamical systems  $\mathbf{z} \sim q_\theta(\mathbf{z})$  and their simulated activity from  $\mathbf{x}(t = 0) = [\frac{\sqrt{2}}{2}, -\frac{\sqrt{2}}{2}]$  colored by log probability. **A.** Each dimension of the simulated trajectories throughout time. **B.** The simulated trajectories in phase space.

1011 The KL divergence can be written in terms of entropy of the variational approximation:

$$KL(q_\theta(\mathbf{z}) \parallel p(\mathbf{z} | \mathbf{x})) = \mathbb{E}_{\mathbf{z} \sim q_\theta} [\log(q_\theta(\mathbf{z}))] - \mathbb{E}_{\mathbf{z} \sim q_\theta} [\log(p(\mathbf{z} | \mathbf{x}))] \quad (33)$$

1012

$$= -H(q_\theta) - \mathbb{E}_{\mathbf{z} \sim q_\theta} [\log(p(\mathbf{x} | \mathbf{z})) + \log(p(\mathbf{z})) - \log(p(\mathbf{x}))] \quad (34)$$

1013 Since the marginal distribution of the data  $p(\mathbf{x})$  (or ‘‘evidence’’) is independent of  $\theta$ , variational  
1014 inference is executed by optimizing the remaining expression. This is usually framed as maximizing  
1015 the evidence lower bound (ELBO)

$$\operatorname{argmin}_{q_\theta \in Q} KL(q_\theta \parallel p(\mathbf{z} | \mathbf{x})) = \operatorname{argmax}_{q_\theta \in Q} H(q_\theta) + \mathbb{E}_{\mathbf{z} \sim q_\theta} [\log(p(\mathbf{x} | \mathbf{z})) + \log(p(\mathbf{z}))]. \quad (35)$$

1016 Now, consider the setting where we have chosen a uniform prior, and stipulate a mean-field gaussian  
1017 likelihood on a chosen statistic of the data  $f(\mathbf{x}; \mathbf{z})$

$$p(\mathbf{x} | \mathbf{z}) = \mathcal{N}(f(\mathbf{x}; \mathbf{z}) | \boldsymbol{\mu}_f, \Sigma_f), \quad (36)$$

1018 where  $\Sigma_f = \text{diag}(\sigma_f^2)$ . The log likelihood is then proportional to a dot product of the natural  
1019 parameter of this mean-field gaussian distribution and the first and second moment statistics.

$$\log p(\mathbf{x} | \mathbf{z}) \propto \boldsymbol{\eta}_f^\top T(\mathbf{x}, \mathbf{z}), \quad (37)$$

1020 where

$$\boldsymbol{\eta}_f = \begin{bmatrix} \frac{\boldsymbol{\mu}_f}{\sigma_f^2} \\ \frac{-1}{2\sigma_f^2} \end{bmatrix}, \text{ and} \quad (38)$$

1021

$$T(\mathbf{x}; \mathbf{z}) = \begin{bmatrix} f(\mathbf{x}; \mathbf{z}) \\ (f(\mathbf{x}; \mathbf{z}) - \boldsymbol{\mu}_f)^2 \end{bmatrix}. \quad (39)$$

1022 The variational objective is then

$$\operatorname{argmax}_{q_{\theta} \in Q} H(q_{\theta}) + \boldsymbol{\eta}_f^{\top} \mathbb{E}_{\mathbf{z} \sim q_{\theta}} [T(\mathbf{x}; \mathbf{z})]. \quad (40)$$

1023 Comparing this to the lagrangian objective (without augmentation) of EPI, we see they are the

1024 same

$$\begin{aligned} q_{\theta}^*(\mathbf{z}) &= \operatorname{argmin}_{q_{\theta} \in Q} -H(q_{\theta}) + \boldsymbol{\eta}_{\text{opt}}^{\top} (\mathbb{E}_{\mathbf{z}, \mathbf{x}} [T(\mathbf{x}; \mathbf{z})] - \boldsymbol{\mu}_{\text{opt}}) \\ &= \operatorname{argmin}_{q_{\theta} \in Q} -H(q_{\theta}) + \boldsymbol{\eta}_{\text{opt}}^{\top} \mathbb{E}_{\mathbf{z}, \mathbf{x}} [T(\mathbf{x}; \mathbf{z})]. \end{aligned} \quad (41)$$

1025 where  $T(\mathbf{x}; \mathbf{z})$  consists of the first and second moments of the emergent property statistic  $f(\mathbf{x}; \mathbf{z})$   
 1026 (Equation 20). Thus, EPI is implicitly executing variational inference with a uniform prior and a  
 1027 mean-field gaussian likelihood on the emergent property statistics. The mean and variances of the  
 1028 mean-field gaussian likelihood are predicated by  $\boldsymbol{\eta}_{\text{opt}}$  (Equations 38 and 40), which is adapted after  
 1029 each EPI optimization epoch based on  $\mathcal{X}$  (see Section 5.1.4). In EPI, the inferred distribution is  
 1030 not conditioned on a finite dataset as in variational inference, but rather the emergent property  
 1031  $\mathcal{X}$  dictates the likelihood parameterization such that the inferred distribution will produce the  
 1032 emergent property. As a note, we could not simply choose  $\boldsymbol{\mu}_f$  and  $\boldsymbol{\sigma}_f$  directly from the outset, since  
 1033 we do not know which of these choices will produce the emergent property  $\mathcal{X}$ , which necessitates  
 1034 the EPI optimization routine that adapts  $\boldsymbol{\eta}_{\text{opt}}$ . Accordingly, we replace the notation of  $p(\mathbf{z} | \mathbf{x})$   
 1035 with  $p(\mathbf{z} | \mathcal{X})$  conceptualizing an inferred distribution that obeys emergent property  $\mathcal{X}$  (see Section  
 1036 5.1).

1037 

## 5.2 Stomatogastric ganglion

1038 In Section 3.1 and 3.2, we used EPI to infer conductance parameters in a model of the stomatogastric  
 1039 ganglion (STG) [41]. This 5-neuron circuit model represents two subcircuits: that generating the  
 1040 pyloric rhythm (fast population) and that generating the gastric mill rhythm (slow population).  
 1041 The additional neuron (the IC neuron of the STG) receives inhibitory synaptic input from both  
 1042 subcircuits, and can couple to either rhythm dependent on modulatory conditions. There is also  
 1043 a parametric regime in which this neuron fires at an intermediate frequency between that of the  
 1044 fast and slow populations [41], which we infer with EPI as a motivational example. This model

1045 is not to be confused with an STG subcircuit model of the pyloric rhythm [69], which has been  
 1046 statistically inferred in other studies [15, 35].

### 1047 5.2.1 STG model

1048 We analyze how the parameters  $\mathbf{z} = [g_{el}, g_{synA}]$  govern the emergent phenomena of intermediate  
 1049 hub frequency in a model of the stomatogastric ganglion (STG) [41] shown in Figure 1A with  
 1050 activity  $\mathbf{x} = [x_{f1}, x_{f2}, x_{hub}, x_{s1}, x_{s2}]$ , using the same hyperparameter choices as Gutierrez et al.  
 1051 Each neuron's membrane potential  $x_\alpha(t)$  for  $\alpha \in \{f1, f2, \text{hub}, s1, s2\}$  is the solution of the following  
 1052 stochastic differential equation:

$$C_m \frac{dx_\alpha}{dt} = -[h_{leak}(\mathbf{x}; \mathbf{z}) + h_{Ca}(\mathbf{x}; \mathbf{z}) + h_K(\mathbf{x}; \mathbf{z}) + h_{hyp}(\mathbf{x}; \mathbf{z}) + h_{elec}(\mathbf{x}; \mathbf{z}) + h_{syn}(\mathbf{x}; \mathbf{z})] + dB. \quad (42)$$

1053 The input current of each neuron is the sum of the leak, calcium, potassium, hyperpolarization,  
 1054 electrical and synaptic currents. Each current component is a function of all membrane potentials  
 1055 and the conductance parameters  $\mathbf{z}$ . Finally, we include gaussian noise  $dB$  to the model of Gutierrez  
 1056 et al. so that the model stochastic, although this is not required by EPI.

1057 The capacitance of the cell membrane was set to  $C_m = 1nF$ . Specifically, the currents are the  
 1058 difference in the neuron's membrane potential and that current type's reversal potential multiplied  
 1059 by a conductance:

$$h_{leak}(\mathbf{x}; \mathbf{z}) = g_{leak}(x_\alpha - V_{leak}) \quad (43)$$

$$h_{elec}(\mathbf{x}; \mathbf{z}) = g_{el}(x_\alpha^{post} - x_\alpha^{pre}) \quad (44)$$

$$h_{syn}(\mathbf{x}; \mathbf{z}) = g_{syn}S_\infty^{pre}(x_\alpha^{post} - V_{syn}) \quad (45)$$

$$h_{Ca}(\mathbf{x}; \mathbf{z}) = g_{Ca}M_\infty(x_\alpha - V_{Ca}) \quad (46)$$

$$h_K(\mathbf{x}; \mathbf{z}) = g_KN(x_\alpha - V_K) \quad (47)$$

$$h_{hyp}(\mathbf{x}; \mathbf{z}) = g_hH(x_\alpha - V_{hyp}). \quad (48)$$

1065 The reversal potentials were set to  $V_{leak} = -40mV$ ,  $V_{Ca} = 100mV$ ,  $V_K = -80mV$ ,  $V_{hyp} = -20mV$ ,  
 1066 and  $V_{syn} = -75mV$ . The other conductance parameters were fixed to  $g_{leak} = 1 \times 10^{-4}\mu S$ .  $g_{Ca}$ ,  
 1067  $g_K$ , and  $g_{hyp}$  had different values based on fast, intermediate (hub) or slow neuron. The fast  
 1068 conductances had values  $g_{Ca} = 1.9 \times 10^{-2}$ ,  $g_K = 3.9 \times 10^{-2}$ , and  $g_{hyp} = 2.5 \times 10^{-2}$ . The intermediate  
 1069 conductances had values  $g_{Ca} = 1.7 \times 10^{-2}$ ,  $g_K = 1.9 \times 10^{-2}$ , and  $g_{hyp} = 8.0 \times 10^{-3}$ . Finally, the  
 1070 slow conductances had values  $g_{Ca} = 8.5 \times 10^{-3}$ ,  $g_K = 1.5 \times 10^{-2}$ , and  $g_{hyp} = 1.0 \times 10^{-2}$ .

1071 Furthermore, the Calcium, Potassium, and hyperpolarization channels have time-dependent gating  
 1072 dynamics dependent on steady-state gating variables  $M_\infty$ ,  $N_\infty$  and  $H_\infty$ , respectively:

$$M_\infty = 0.5 \left( 1 + \tanh \left( \frac{x_\alpha - v_1}{v_2} \right) \right) \quad (49)$$

$$\frac{dN}{dt} = \lambda_N (N_\infty - N) \quad (50)$$

$$N_\infty = 0.5 \left( 1 + \tanh \left( \frac{x_\alpha - v_3}{v_4} \right) \right) \quad (51)$$

$$\lambda_N = \phi_N \cosh \left( \frac{x_\alpha - v_3}{2v_4} \right) \quad (52)$$

$$\frac{dH}{dt} = \frac{(H_\infty - H)}{\tau_h} \quad (53)$$

$$H_\infty = \frac{1}{1 + \exp \left( \frac{x_\alpha + v_5}{v_6} \right)} \quad (54)$$

$$\tau_h = 272 - \left( \frac{-1499}{1 + \exp \left( \frac{-x_\alpha + v_7}{v_8} \right)} \right). \quad (55)$$

1079 where we set  $v_1 = 0mV$ ,  $v_2 = 20mV$ ,  $v_3 = 0mV$ ,  $v_4 = 15mV$ ,  $v_5 = 78.3mV$ ,  $v_6 = 10.5mV$ ,  
 1080  $v_7 = -42.2mV$ ,  $v_8 = 87.3mV$ ,  $v_9 = 5mV$ , and  $v_{th} = -25mV$ .

1081 Finally, there is a synaptic gating variable as well:

$$S_\infty = \frac{1}{1 + \exp \left( \frac{v_{th} - x_\alpha}{v_9} \right)}. \quad (56)$$

1082 When the dynamic gating variables are considered, this is actually a 15-dimensional nonlinear  
 1083 dynamical system. The gaussian noise  $d\mathbf{B}$  has variance  $(1 \times 10^{-12})^2$  A<sup>2</sup>, and introduces variability  
 1084 in frequency at each parameterization  $\mathbf{z}$ .

### 1085 5.2.2 Hub frequency calculation

1086 In order to measure the frequency of the hub neuron during EPI, the STG model was simulated for  
 1087  $T = 300$  time steps of  $dt = 25\text{ms}$ . The chosen  $dt$  and  $T$  were the most computationally convenient  
 1088 choices yielding accurate frequency measurement. We used a basis of complex exponentials with  
 1089 frequencies from 0.0-1.0 Hz at 0.01Hz resolution to measure frequency from simulated time series

$$\Phi = [0.0, 0.01, \dots, 1.0]^\top .. \quad (57)$$

1090 To measure spiking frequency, we processed simulated membrane potentials with a relu (spike  
 1091 extraction) and low-pass filter with averaging window of size 20, then took the frequency with the

1092 maximum absolute value of the complex exponential basis coefficients of the processed time-series.  
 1093 The first 20 temporal samples of the simulation are ignored to account for initial transients.

1094 To differentiate through the maximum frequency identification, we used a soft-argmax Let  $X_\alpha \in$   
 1095  $\mathcal{C}^{|\Phi|}$  be the complex exponential filter bank dot products with the signal  $x_\alpha \in \mathbb{R}^N$ , where  $\alpha \in$   
 1096  $\{f1, f2, \text{hub}, s1, s2\}$ . The soft-argmax is then calculated using temperature parameter  $\beta_\psi = 100$

$$\psi_\alpha = \text{softmax}(\beta_\psi |X_\alpha| \odot i), \quad (58)$$

1097 where  $i = [0, 1, \dots, 100]$ . The frequency is then calculated as

$$\omega_\alpha = 0.01\psi_\alpha \text{Hz}. \quad (59)$$

1098 Intermediate hub frequency, like all other emergent properties in this work, is defined by the mean  
 1099 and variance of the emergent property statistics. In this case, we have one statistic, hub neuron  
 1100 frequency, where the mean was chosen to be 0.55Hz,(Equation 2) and variance was chosen to be  
 1101  $0.025^2$  Hz $^2$  (Equation 3).

1102 **5.2.3 EPI details for the STG model**

1103 As a maximum entropy distribution,  $T(\mathbf{x}; \mathbf{z})$  is comprised of both these first and second moments  
 1104 of the hub neuron frequency (as in Equations 20 and 21)

$$T(\mathbf{x}; \mathbf{z}) = \begin{bmatrix} \omega_{\text{hub}}(\mathbf{x}; \mathbf{z}) \\ (\omega_{\text{hub}}(\mathbf{x}; \mathbf{z}) - 0.55)^2 \end{bmatrix}, \quad (60)$$

$$\boldsymbol{\mu}_{\text{opt}} = \begin{bmatrix} 0.55 \\ 0.025^2 \end{bmatrix}. \quad (61)$$

1105 1106 Throughout optimization, the augmented lagrangian parameters  $\eta$  and  $c$ , were updated after each  
 1107 epoch of 5,000 iterations(see Section 5.1.4). The optimization converged after five epochs (Fig. S4).

1108 For EPI in Fig 1E, we used a real NVP architecture with three Real NVP coupling layers and two-  
 1109 layer neural networks of 25 units per layer. The normalizing flow architecture mapped  $\mathbf{z}_0 \sim \mathcal{N}(\mathbf{0}, I)$   
 1110 to a support of  $\mathbf{z} = [g_{\text{el}}, g_{\text{synA}}] \in [4, 8] \times [0.01, 4]$ , initialized to a gaussian approximation of samples  
 1111 returned by a preliminary ABC search. We did not include  $g_{\text{synA}} < 0.01$ , for numerical stability.  
 1112 EPI optimization was run using 5 different random seeds for architecture initialization  $\boldsymbol{\theta}$  with an  
 1113 augmented lagrangian coefficient of  $c_0 = 10^5$ , a batch size  $n = 400$ , and  $\beta = 2$ . The architecture  
 1114 converged with criteria  $N_{\text{test}} = 100$ .

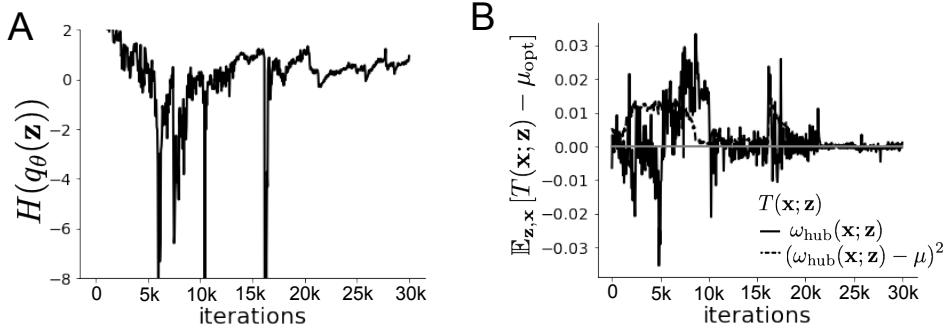


Figure S4: EPI optimization of the STG model producing network syncing. **A.** Entropy throughout optimization. **B.** The emergent property statistic means and variances converge to their constraints at 25,000 iterations following the fifth augmented lagrangian epoch.

1115 **5.2.4 Hessian sensitivity vectors**

1116 To quantify the second-order structure of the EPI distribution, we evaluated the Hessian of the log  
1117 probability  $\frac{\partial^2 \log q(\mathbf{z}|\mathcal{X})}{\partial \mathbf{z} \partial \mathbf{z}^\top}$ . The eigenvector of this Hessian with most negative eigenvalue is defined as  
1118 the sensitivity dimension  $\mathbf{v}_1$ , and all subsequent eigenvectors are ordered by increasing eigenvalue.  
1119 These eigenvalues are quantifications of how fast the emergent property deteriorates via the param-  
1120 eter combination of their associated eigenvector. In Figure 1D, the sensitivity dimension  $v_1$  (solid)  
1121 and the second eigenvector of the Hessian  $v_2$  (dashed) are shown evaluated at the mode of the  
1122 distribution. Since the Hessian eigenvectors have sign degeneracy, the visualized directions in 2-D  
1123 parameter space were chosen to have positive  $g_{\text{synA}}$ . The length of the arrows is inversely propor-  
1124 tional to the square root of the absolute value of their eigenvalues  $\lambda_1 = -10.7$  and  $\lambda_2 = -3.22$ . For  
1125 the same magnitude perturbation away from the mode, intermediate hub frequency only diminishes  
1126 along the sensitivity dimension  $\mathbf{v}_1$  (Fig. 1E-F).

1127 **5.3 Scaling EPI for stable amplification in RNNs**

1128 **5.3.1 Rank-2 RNN model**

1129 We examined the scaling properties of EPI by learning connectivities of RNNs of increasing size  
1130 that exhibit stable amplification. Rank-2 RNN connectivity was modeled as  $W = UV^\top$ , where  
1131  $U = [\mathbf{U}_1 \quad \mathbf{U}_2] + g\chi^{(W)}$ ,  $V = [\mathbf{V}_1 \quad \mathbf{V}_2] + g\chi^{(V)}$ , and  $\chi_{i,j}^{(W)}, \chi_{i,j}^{(V)} \sim \mathcal{N}(0, 1)$ . This RNN model has  
1132 dynamics

$$\tau \dot{\mathbf{x}} = -\mathbf{x} + W\mathbf{x}. \quad (62)$$

1133 In this analysis, we inferred connectivity parameterizations  $\mathbf{z} = [\mathbf{U}_1^\top, \mathbf{U}_2^\top, \mathbf{V}_1^\top, \mathbf{V}_2^\top]^\top \in [-1, 1]^{(4N)}$   
1134 that produced stable amplification using EPI, SMC-ABC [26], and SNPE [35] (see Section Related  
1135 Methods).

1136 **5.3.2 Stable amplification**

1137 For this RNN model to be stable, all real eigenvalues of  $W$  must be less than 1:  $\text{real}(\lambda_1) < 1$ ,  
1138 where  $\lambda_1$  denotes the greatest real eigenvalue of  $W$ . For a stable RNN to amplify at least one input  
1139 pattern, the symmetric connectivity  $W^s = \frac{W + W^\top}{2}$  must have an eigenvalue greater than 1:  $\lambda_1^s > 1$ ,  
1140 where  $\lambda^s$  is the maximum eigenvalue of  $W^s$ . These two conditions are necessary and sufficient for  
1141 stable amplification in RNNs [51].

1142 **5.3.3 EPI details for RNNs**

1143 We defined the emergent property of stable amplification with means of these eigenvalues (0.5  
1144 and 1.5, respectively) that satisfy these conditions. To complete the emergent property definition,  
1145 we chose variances ( $0.25^2$ ) about those means such that samples rarely violate the eigenvalue  
1146 constraints. In terms of the EPI optimization variables, this is written as

$$T(\mathbf{x}; \mathbf{z}) = \begin{bmatrix} \text{real}(\lambda_1)(\mathbf{x}; \mathbf{z}) \\ \lambda_1^s(\mathbf{x}; \mathbf{z}) \\ (\text{real}(\lambda_1)(\mathbf{x}; \mathbf{z}) - 0.5)^2 \\ (\lambda_1^s(\mathbf{x}; \mathbf{z}) - 1.5)^2 \end{bmatrix}, \quad (63)$$

1147

$$\boldsymbol{\mu}_{\text{opt}} = \begin{bmatrix} 0.5 \\ 1.5 \\ 0.25^2 \\ 0.25^2 \end{bmatrix}. \quad (64)$$

1148 Gradients of maximum eigenvalues of Hermitian matrices like  $W^s$  are available with modern auto-  
1149 matic differentiation tools. To differentiate through the  $\text{real}(\lambda_1)$ , we solved the following equation  
1150 for eigenvalues of rank-2 matrices using the rank reduced matrix  $W^r = V^\top U$

$$\lambda_{\pm} = \frac{\text{Tr}(W^r) \pm \sqrt{\text{Tr}(W^r)^2 - 4\text{Det}(W^r)}}{2}. \quad (65)$$

1151 For EPI in Fig. 2, we used a real NVP architecture with three coupling layers of affine transfor-  
1152 mations parameterized by two-layer neural networks of 100 units per layer. The initial distribution

1153 was a standard isotropic gaussian  $\mathbf{z}_0 \sim \mathcal{N}(\mathbf{0}, I)$  mapped to the support of  $\mathbf{z}_i \in [-1, 1]$ . We used  
 1154 an augmented lagrangian coefficient of  $c_0 = 10^3$ , a batch size  $n = 200$ ,  $\beta = 4$ , and chose to use  
 1155 500 iterations per augmented lagrangian epoch and emergent property constraint convergence was  
 1156 evaluated at  $N_{\text{test}} = 200$  (Fig. 2B blue line, and Fig. 2C-D blue).

1157 **5.3.4 Methodological comparison**

1158 We compared EPI to two alternative simulation-based inference techniques, since the likelihood  
 1159 of these eigenvalues given  $\mathbf{z}$  is not available. Approximate bayesian computation (ABC) [24] is a  
 1160 rejection sampling technique for obtaining sets of parameters  $\mathbf{z}$  that produce activity  $\mathbf{x}$  close to some  
 1161 observed data  $\mathbf{x}_0$ . Sequential Monte Carlo approximate bayesian computation (SMC-ABC) is the  
 1162 state-of-the-art ABC method, which leverages SMC techniques to improve sampling speed. We ran  
 1163 SMC-ABC with the pyABC package [93] to infer RNNs with stable amplification: connectivities  
 1164 having eigenvalues within an  $\epsilon$ -defined  $l_2$  distance of

$$\mathbf{x}_0 = \begin{bmatrix} \text{real}(\lambda_1) \\ \lambda_1^s \end{bmatrix} = \begin{bmatrix} 0.5 \\ 1.5 \end{bmatrix}. \quad (66)$$

1165 SMC-ABC was run with a uniform prior over  $\mathbf{z} \in [-1, 1]^{(4N)}$ , a population size of 1,000 particles  
 1166 with simulations parallelized over 32 cores, and a multivariate normal transition model.

1167 SNPE, the next approach in our comparison, is far more similar to EPI. Like EPI, SNPE treats pa-  
 1168 rameters in mechanistic models with deep probability distributions, yet the two learning algorithms  
 1169 are categorically different. SNPE uses a two-network architecture to approximate the posterior dis-  
 1170 tribution of the model conditioned on observed data  $\mathbf{x}_0$ . The amortizing network maps observations  
 1171  $\mathbf{x}_i$  to the parameters of the deep probability distribution. The weights and biases of the parameter  
 1172 network are optimized by sequentially augmenting the training data with additional pairs  $(\mathbf{z}_i, \mathbf{x}_i)$   
 1173 based on the most recent posterior approximation. This sequential procedure is important to get  
 1174 training data  $\mathbf{z}_i$  to be closer to the true posterior, and  $\mathbf{x}_i$  to be closer to the observed data. For  
 1175 the deep probability distribution architecture, we chose a masked autoregressive flow with affine  
 1176 couplings (the default choice), three transforms, 50 hidden units, and a normalizing flow mapping  
 1177 to the support as in EPI. This architectural choice closely tracked the size of the architecture used  
 1178 by EPI (Fig. S5). As in SMC-ABC, we ran SNPE with  $\mathbf{x}_0 = \mu$ . All SNPE optimizations were  
 1179 run for a limit of 1.5 days on a Tesla V100 GPU, or until two consecutive rounds resulted in a  
 1180 validation log probability lower than the maximum observed for that random seed.

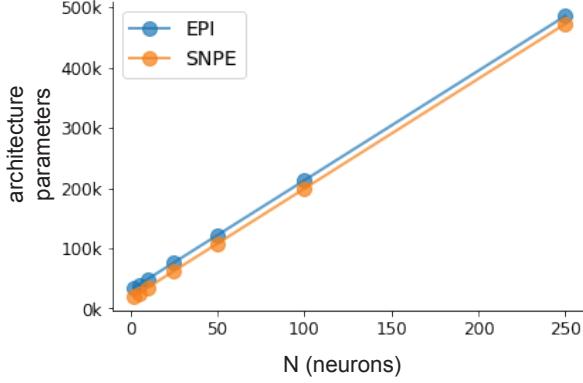


Figure S5: Number of parameters in deep probability distribution architectures of EPI (blue) and SNPE (orange) by RNN size ( $N$ ).

1181 To compare the efficiency of these algorithms for inferring RNN connectivity distributions producing  
 1182 stable amplification, we develop a convergence criteria that can be used across methods. While EPI  
 1183 has its own hypothesis testing convergence criteria for the emergent property, it would not make  
 1184 sense to use this criteria on SNPE and SMC-ABC which do not constrain the means and variances  
 1185 of their predictions. Instead, we consider EPI and SNPE to have converged after completing its  
 1186 most recent optimization epoch (EPI) or round (SNPE) in which the distance

$$d(q_{\theta}(\mathbf{z})) = \|\mathbb{E}_{\mathbf{z}, \mathbf{x}} [f(\mathbf{x}; \mathbf{z})] - \boldsymbol{\mu}\|_2 \quad (67)$$

1187 is less than 0.5. We consider SMC-ABC to have converged once the population produces samples  
 1188 within the  $\epsilon = 0.5$  ball ensuring stable amplification.

1189 When assessing the scalability of SNPE, it is important to check that alternative hyperparameter-  
 1190izations could not yield better performance. Key hyperparameters of the SNPE optimization are  
 1191 the number of simulations per round  $n_{\text{round}}$ , the number of atoms used in the atomic proposals of  
 1192 the SNPE-C algorithm [94], and the batch size  $n$ . To match EPI, we used a batch size of  $n = 200$   
 1193 for  $N \leq 25$ , however we found  $n = 1,000$  to be helpful for SNPE in higher dimensions. While  
 1194  $n_{\text{round}} = 1,000$  yielded SNPE convergence for  $N \leq 25$ , we found that a substantial increase to  
 1195  $n_{\text{round}} = 25,000$  yielded more consistent convergence at  $N = 50$  (Fig. S6A). By increasing  $n_{\text{round}}$ ,  
 1196 we also necessarily increase the duration of each round. At  $N = 100$ , we tried two hyperparameter  
 1197 modifications. As suggested in [94], we increased  $n_{\text{atom}}$  by an order of magnitude to improve gra-  
 1198 dient quality, but this had little effect on the optimization (much overlap between same random  
 1199 seeds) (Fig. S6B). Finally, we increased  $n_{\text{round}}$  by an order of magnitude, which yielded conver-  
 1200 gence in one case, but no others. We found no way to improve the convergence rate of SNPE

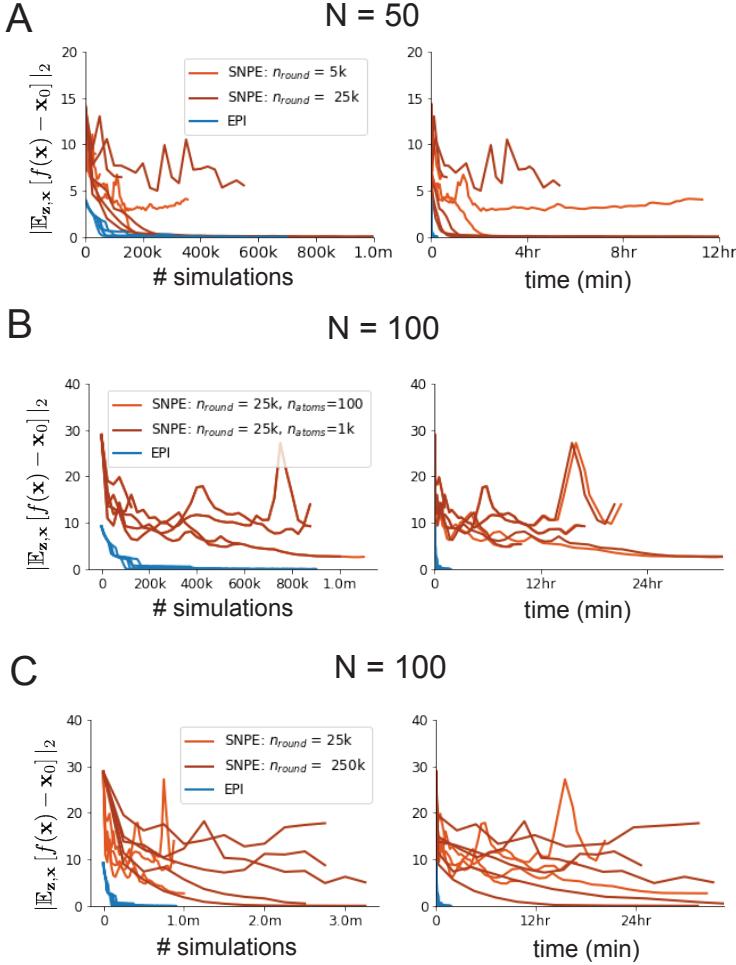


Figure S6: SNPE convergence was enabled by increasing  $n_{\text{round}}$ , not  $n_{\text{atom}}$ . **A.** Difference of mean predictions  $\mathbf{x}_0$  throughout optimization at  $N = 50$  with by simulation count (left) and wall time (right) of SNPE with  $n_{\text{round}} = 5,000$  (light orange), SNPE with  $n_{\text{round}} = 25,000$  (dark orange), and EPI (blue). Each line shows an individual random seed. **B.** Same conventions as A at  $N = 100$  of SNPE with  $n_{\text{atom}} = 100$  (light orange) and  $n_{\text{atom}} = 1,000$  (dark orange). **C.** Same conventions as A at  $N = 100$  of SNPE with  $n_{\text{round}} = 25,000$  (light orange) and  $n_{\text{round}} = 250,000$  (dark orange).

1201 without making more aggressive hyperparameter choices requiring high numbers of simulations. In  
1202 Figure 2C-D, we show samples from the random seed resulting in emergent property convergence  
1203 at greatest entropy (EPI), the random seed resulting in greatest validation log probability (SNPE),  
1204 and the result of all converged random seeds (SMC).

1205 **5.3.5 Effect of RNN parameters on EPI and SNPE inferred distributions**

1206 To clarify the difference in objectives of EPI and SNPE, we show their results on RNN models  
1207 with different numbers of neurons  $N$  and random strength  $g$ . The parameters inferred by EPI  
1208 consistently produces the same mean and variance of  $\text{real}(\lambda_1)$  and  $\lambda_1^s$ , while those inferred by  
1209 SNPE change according to the model definition (Fig. S7A). For  $N = 2$  and  $g = 0.01$ , the SNPE  
1210 posterior has greater concentration in eigenvalues around  $\mathbf{x}_0$  than at  $g = 0.1$ , where the model has  
1211 greater randomness (Fig. S7B top, orange). At both levels of  $g$  when  $N = 2$ , the posterior of SNPE  
1212 has lower entropy than EPI at convergence (Fig. S7B top). However at  $N = 10$ , SNPE results in  
1213 a predictive distribution of more widely dispersed eigenvalues (Fig. S7A bottom), and an inferred  
1214 posterior with greater entropy than EPI (Fig. S7B bottom). We highlight these differences not  
1215 to focus on an insightful trend, but to emphasize that these methods optimize different objectives  
1216 with different implications.

1217 Note that SNPE converges when its validation log probability has saturated after several rounds  
1218 of optimization (Fig. S7C), and that EPI converges after several epochs of its own optimization  
1219 to enforce the emergent property constraints (Fig. S7D blue). Importantly, as SNPE optimizes  
1220 its posterior approximation, the predictive means change, and at convergence may be different  
1221 than  $\mathbf{x}_0$  (Fig. S7D orange, left). It is sensible to assume that predictions of a well-approximated  
1222 SNPE posterior should closely reflect the data on average (especially given a uniform prior and  
1223 a low degree of stochasticity), however this is not a given. Furthermore, no aspect of the SNPE  
1224 optimization controls the variance of the predictions (Fig. S7D orange, right).

1225 **5.4 Primary visual cortex**

1226 **5.4.1 V1 model**

1227 E-I circuit models, rely on the assumption that inhibition can be studied as an indivisible unit,  
1228 despite ample experimental evidence showing that inhibition is instead composed of distinct ele-  
1229 ments [63]. In particular three types of genetically identified inhibitory cell-types – parvalbumin

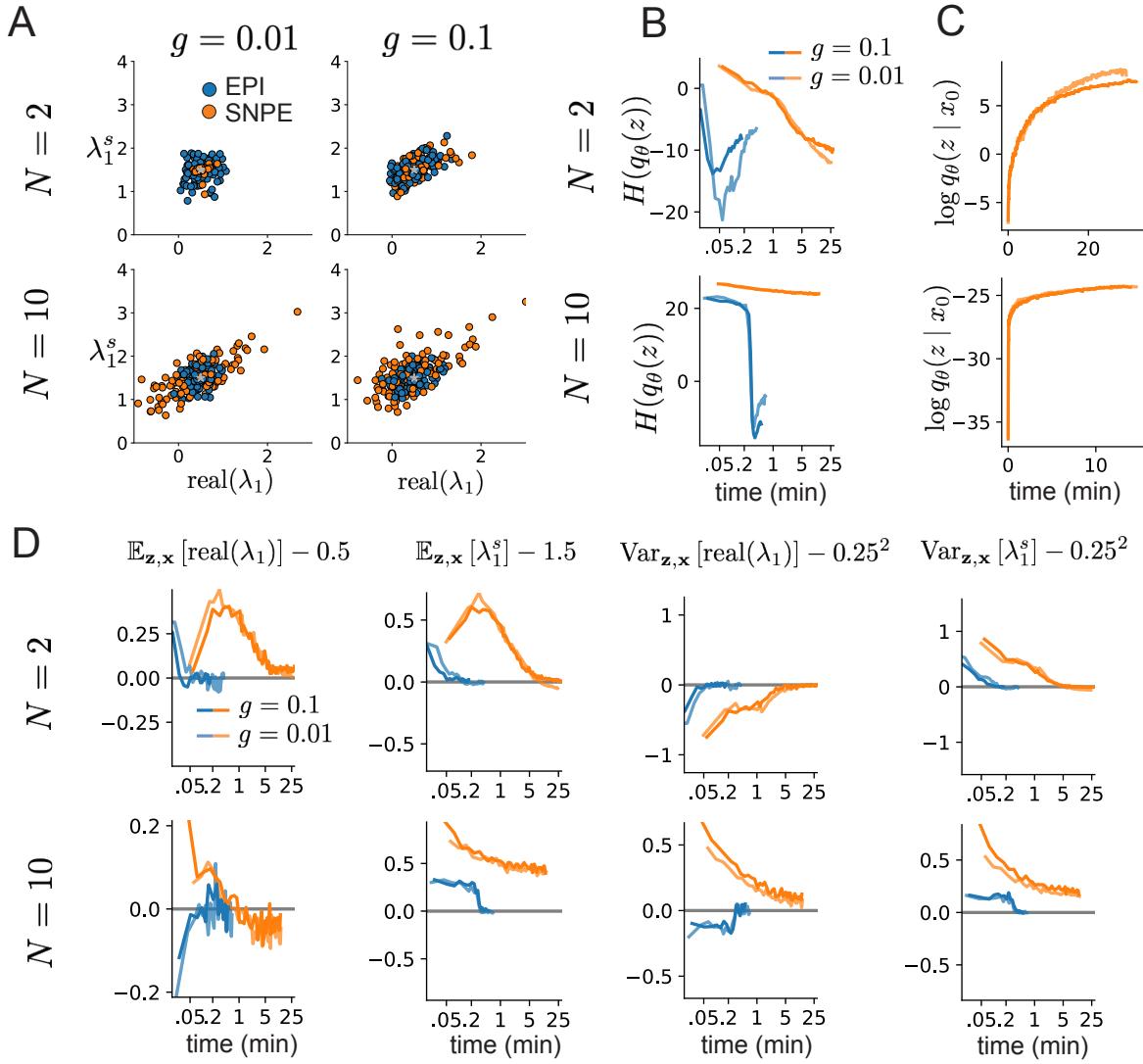


Figure S7: Model characteristics affect predictions of posteriors inferred by SNPE, while predictions of parameters inferred by EPI remain fixed. **A.** Predictive distribution of EPI (blue) and SNPE (orange) inferred connectivity of RNNs exhibiting stable amplification with  $N = 2$  (top),  $N = 10$  (bottom),  $g = 0.01$  (left), and  $g = 0.1$  (right). **B.** Entropy of parameter distribution approximations throughout optimization with  $N = 2$  (top),  $N = 10$  (bottom),  $g = 0.1$  (dark shade), and  $g = 0.01$  (light shade). **C.** Validation log probabilities throughout SNPE optimization. Same conventions as B. **D.** Adherence to EPI constraints. Same conventions as B.

1230 (P), somatostatin (S), VIP (V) – compose 80% of GABAergic interneurons in V1 [61–63], and follow  
 1231 specific connectivity patterns (Fig. 3A) [64], which lead to cell-type specific computations [47, 95].  
 1232 Currently, how the subdivision of inhibitory cell-types, shapes correlated variability by reconfigur-  
 1233 ing recurrent network dynamics is not understood.

1234 In the stochastic stabilized supralinear network [59], population rate responses  $\mathbf{x}$  to mean input  $\mathbf{h}$ ,  
 1235 recurrent input  $W\mathbf{x}$  and slow noise  $\boldsymbol{\epsilon}$  are governed by

$$\tau \frac{d\mathbf{x}}{dt} = -\mathbf{x} + \phi(W\mathbf{x} + \mathbf{h} + \boldsymbol{\epsilon}), \quad (68)$$

1236 where the noise is an Ornstein-Uhlenbeck process  $\boldsymbol{\epsilon} \sim OU(\tau_{\text{noise}}, \boldsymbol{\sigma})$

$$\tau_{\text{noise}} d\epsilon_\alpha = -\epsilon_\alpha dt + \sqrt{2\tau_{\text{noise}}} \tilde{\sigma}_\alpha dB \quad (69)$$

1237 with  $\tau_{\text{noise}} = 5\text{ms} > \tau = 1\text{ms}$ . The noisy process is parameterized as

$$\tilde{\sigma}_\alpha = \sigma_\alpha \sqrt{1 + \frac{\tau}{\tau_{\text{noise}}}}, \quad (70)$$

1238 so that  $\boldsymbol{\sigma}$  parameterizes the variance of the noisy input in the absence of recurrent connectivity  
 1239 ( $W = \mathbf{0}$ ). As contrast  $c \in [0, 1]$  increases, input to the E- and P-populations increases relative to  
 1240 a baseline input  $\mathbf{h} = \mathbf{h}_b + c\mathbf{h}_c$ . Connectivity ( $W_{\text{fit}}$ ) and input ( $\mathbf{h}_{b,\text{fit}}$  and  $\mathbf{h}_{c,\text{fit}}$ ) parameters were fit  
 1241 using the deterministic V1 circuit model [47]

$$W_{\text{fit}} = \begin{bmatrix} W_{EE} & W_{EP} & W_{ES} & W_{EV} \\ W_{PE} & W_{PP} & W_{PS} & W_{PV} \\ W_{SE} & W_{SP} & W_{SS} & W_{SV} \\ W_{VE} & W_{VP} & W_{VS} & W_{VV} \end{bmatrix} = \begin{bmatrix} 2.18 & -1.19 & -.594 & -.229 \\ 1.66 & -.651 & -.680 & -.242 \\ .895 & -5.22 \times 10^{-3} & -1.51 \times 10^{-4} & -.761 \\ 3.34 & -2.31 & -.254 & -2.52 \times 10^{-4} \end{bmatrix}, \quad (71)$$

$$\mathbf{h}_{b,\text{fit}} = \begin{bmatrix} .416 \\ .429 \\ .491 \\ .486 \end{bmatrix}, \quad (72)$$

1242 and

$$\mathbf{h}_{c,\text{fit}} = \begin{bmatrix} .359 \\ .403 \\ 0 \\ 0 \end{bmatrix}. \quad (73)$$

1243 To obtain rates on a realistic scale (100-fold greater), we map these fitted parameters to an equivalence class  
1244

$$W = \begin{bmatrix} W_{EE} & W_{EP} & W_{ES} & W_{EV} \\ W_{PE} & W_{PP} & W_{PS} & W_{PV} \\ W_{SE} & W_{SP} & W_{SS} & W_{SV} \\ W_{VE} & W_{VP} & W_{VS} & W_{VV} \end{bmatrix} = \begin{bmatrix} .218 & -.119 & -.0594 & -.0229 \\ .166 & -.0651 & -.068 & -.0242 \\ .0895 & -5.22 \times 10^{-4} & -1.51 \times 10^{-5} & -.0761 \\ .334 & -.231 & -.0254 & -2.52 \times 10^{-5} \end{bmatrix}, \quad (74)$$

$$\mathbf{h}_b = \begin{bmatrix} h_{b,E} \\ h_{b,P} \\ h_{b,S} \\ h_{b,V} \end{bmatrix} = \begin{bmatrix} 4.16 \\ 4.29 \\ 4.91 \\ 4.86 \end{bmatrix}, \quad (75)$$

1245 and

$$\mathbf{h}_c = \begin{bmatrix} h_{c,E} \\ h_{c,P} \\ h_{c,S} \\ h_{c,V} \end{bmatrix} = \begin{bmatrix} 3.59 \\ 4.03 \\ 0 \\ 0 \end{bmatrix}. \quad (76)$$

1246 Circuit responses are simulated using  $T = 200$  time steps at  $dt = 0.5\text{ms}$  from an initial condition  
1247 drawn from  $\mathbf{x}(0) \sim U[10 \text{ Hz}, 25 \text{ Hz}]$ . Standard deviation of the E-population  $s_E(\mathbf{x}; \mathbf{z})$  is calculated  
1248 as the square root of the temporal variance from  $t_{ss} = 75\text{ms}$  to  $Tdt = 100\text{ms}$  averaged over 100  
1249 independent trials.

$$s_E(\mathbf{x}; \mathbf{z}) = \mathbb{E}_x \left[ \sqrt{\mathbb{E}_{t>t_{ss}} [(x_E(t) - \mathbb{E}_{t>t_{ss}} [x_E(t)])^2]} \right] \quad (77)$$

#### 1250 5.4.2 EPI details for the V1 model

1251 For EPI in Figures 3D-E and S8, we used a real NVP architecture with three Real NVP coupling  
1252 layers and two-layer neural networks of 50 units per layer. The normalizing flow architecture  
1253 mapped  $z_0 \sim \mathcal{N}(\mathbf{0}, I)$  to a support of  $\mathbf{z} = [\sigma_E, \sigma_P, \sigma_S, \sigma_V] \in [0.0, 0.5]^4$ . EPI optimization was run  
1254 using three different random seeds for architecture initialization  $\boldsymbol{\theta}$  with an augmented lagrangian  
1255 coefficient of  $c_0 = 10^{-1}$ , a batch size  $n = 100$ , and  $\beta = 2$ . The distributions shown are those of the  
1256 architectures converging with criteria  $N_{\text{test}} = 100$  at greatest entropy across three random seeds.

1257 Optimization details are shown in Figure S9. The sums of squares of each pair of parameters are  
 1258 shown for each EPI distribution in Figure S10.

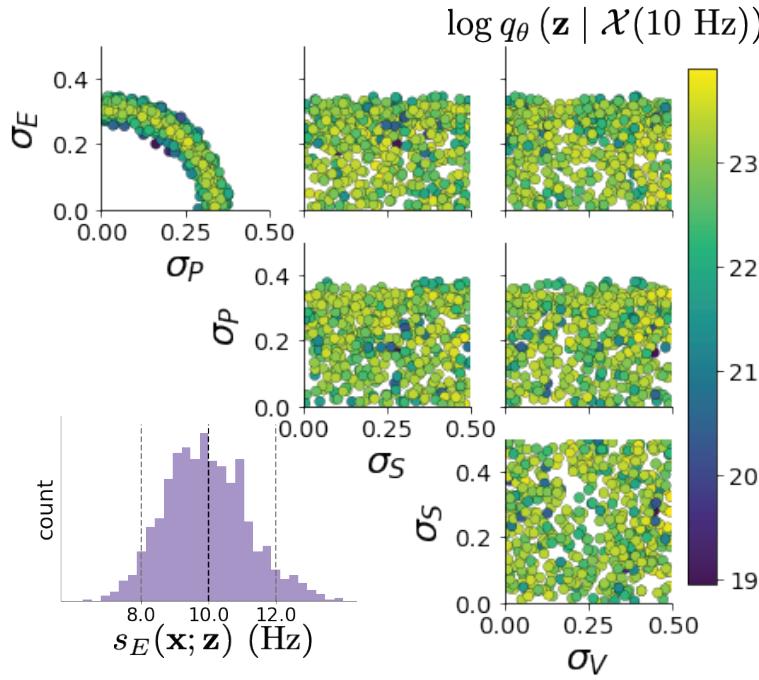


Figure S8: EPI inferred distribution for  $\mathcal{X}(10 \text{ Hz})$ .

### 1259 5.4.3 Sensitivity analyses

1260 In Fig. 3E, we visualize the modes of  $q_{\theta}(\mathbf{z} \mid \mathcal{X})$  throughout the  $\sigma_E$ - $\sigma_P$  marginal. Specifically, we  
 1261 calculated

$$\begin{aligned} \mathbf{z}^*(\sigma_{P,\text{fixed}}) &= \underset{\mathbf{z}}{\operatorname{argmax}} \log q_{\theta}(\mathbf{z} \mid \mathcal{X}) \\ \text{s.t. } \sigma_P &= \sigma_{P,\text{fixed}} \end{aligned} \quad (78)$$

1262 At each mode  $\mathbf{z}^*$ , we calculated the Hessian and visualized the sensitivity dimension in the direction  
 1263 of positive  $\sigma_E$ .

### 1264 5.4.4 Testing for the paradoxical effect

1265 The paradoxical effect occurs when a populations steady state rate is decreased (or increased)  
 1266 when an increase (decrease) in current is applied to that population [12]. To see which, if any,  
 1267 populations exhibited a paradoxical effect, we examined responses to changes in input (Fig. S11).

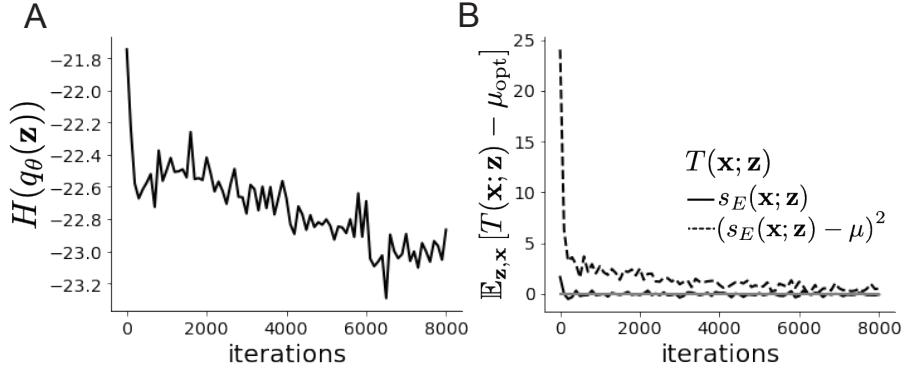


Figure S9: EPI optimization  $q_\theta(\mathbf{z} | \mathcal{X}(5 \text{ Hz}))$  **A.** Entropy throughout optimization. **B.** The emergent property statistic means and variances converge to their constraints at 8,000 iterations following the fourth augmented lagrangian epoch.

1268 Input magnitudes were chosen so that the effect is salient (0.002 for E and P, but 0.02 for S and  
1269 V). Only the P-population exhibited the paradoxical effect at this connectivity  $W$  and input  $\mathbf{h}$ .

1270 **5.4.5 Primary visual cortex: Mathematical intuition and challenges**

1271 The dynamical system that we are working with can be written as

$$\begin{aligned} dx &= \frac{1}{\tau}(-x + f(Wx + h + \epsilon))dt \\ d\epsilon &= -\frac{dt}{\tau_{\text{noise}}} \epsilon + \frac{\sqrt{2}}{\sqrt{\tau_{\text{noise}}}} \Sigma_\epsilon dW \end{aligned} \tag{79}$$

1272 Where in this paper we chose

$$\Sigma_\epsilon = \tau_{\text{noise}} \begin{bmatrix} \tilde{\sigma}_E & 0 & 0 & 0 \\ 0 & \tilde{\sigma}_P & 0 & 0 \\ 0 & 0 & \tilde{\sigma}_S & 0 \\ 0 & 0 & 0 & \tilde{\sigma}_V \end{bmatrix} \tag{80}$$

1273 where  $\tilde{\sigma}_\alpha$  is the reparameterized standard deviation of the noise for population  $\alpha$  from Equation  
1274 70.

1275 In order to compute this covariance, we define  $v = \omega x + h + \epsilon$  and  $S = I - \omega f'(v)$ , to re-write Eq.  
1276 (79) as an 8-dimensional system:

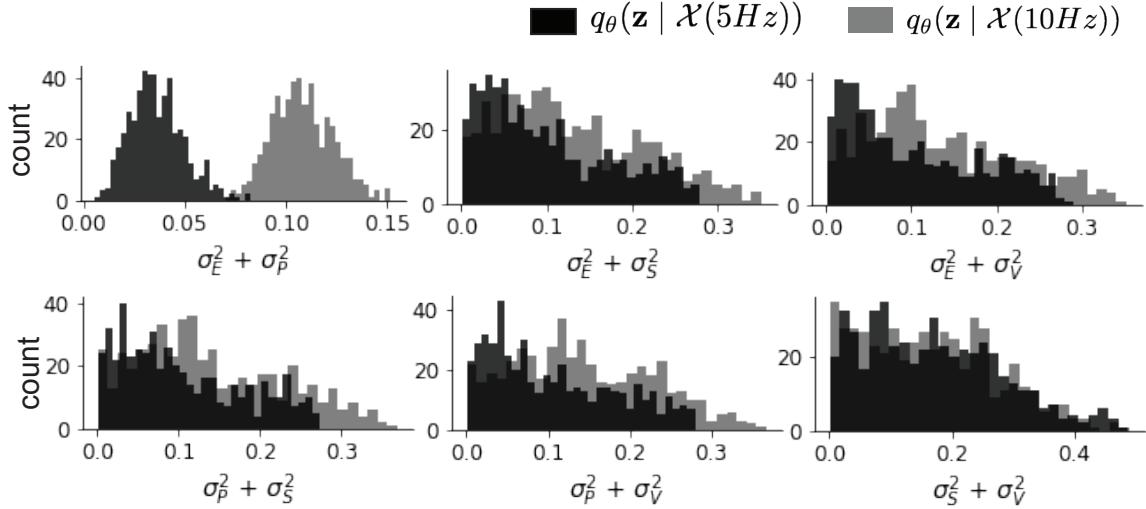


Figure S10: EPI predictive distributions of the sum of squares of each pair of noise parameters.

$$d \begin{pmatrix} \delta v \\ \epsilon \end{pmatrix} = - \begin{pmatrix} S & -\frac{\tau_{\text{noise}} - \tau}{\tau \tau_{\text{noise}}} I \\ 0 & \frac{1}{\tau_{\text{noise}}} I \end{pmatrix} \begin{pmatrix} \delta v \\ \epsilon \end{pmatrix} dt + \begin{pmatrix} 0 & \frac{\sqrt{2}}{\sqrt{\tau_{\text{noise}}}} \Sigma_\epsilon \\ 0 & \frac{\sqrt{2}}{\sqrt{\tau_{\text{noise}}}} \Sigma_\epsilon \end{pmatrix} d\mathbf{W} \quad (81)$$

1277 Where  $d\mathbf{W}$  is a vector with the private noise of each variable. The  $d\mathbf{W}$  term is multiplied by a  
 1278 non-diagonal matrix is because the noise that the voltage receives is the exact same than the one  
 1279 that comes from the OU process and not another process. The solution of this problem is given by  
 1280 the Lyapunov Equation [59, 66]:

$$\begin{pmatrix} S & -\frac{\tau_{\text{noise}} - \tau}{\tau \tau_{\text{noise}}} I \\ 0 & \frac{1}{\tau_{\text{noise}}} I \end{pmatrix} \begin{pmatrix} \Lambda_v & \Lambda_c \\ \Lambda_c^T & \Lambda_\epsilon \end{pmatrix} + \begin{pmatrix} \Lambda_v & \Lambda_c \\ \Lambda_c^T & \Lambda_\epsilon \end{pmatrix} \begin{pmatrix} S^T & 0 \\ -\frac{\tau_{\text{noise}} - \tau}{\tau \tau_{\text{noise}}} I & \frac{1}{\tau_{\text{noise}}} I \end{pmatrix} = \begin{pmatrix} \frac{2}{\tau_{\text{noise}}} \Lambda_\epsilon & \frac{2}{\tau_{\text{noise}}} \Lambda_\epsilon \\ \frac{2}{\tau_{\text{noise}}} \Lambda_\epsilon & \frac{2}{\tau_{\text{noise}}} \Lambda_\epsilon \end{pmatrix} \quad (82)$$

1281 To obtain an equation for  $\Lambda_v$ , we solve this block matrix multiplication:

$$S\Lambda_v + \Lambda_v S^T = \frac{2\Lambda_\epsilon}{\tau_{\text{noise}}} + \frac{\tau_{\text{noise}}^2 - \tau^2}{(\tau \tau_{\text{noise}})^2} \left( \left( \frac{1}{\tau_{\text{noise}}} I + S \right)^{-1} \Lambda_\epsilon + \Lambda_\epsilon \left( \frac{1}{\tau_{\text{noise}}} I + S^T \right)^{-1} \right) \quad (83)$$

Which is another Lyapunov Equation, now in 4 dimensions. In the simplest case in which  $\tau_{\text{noise}} = \tau$ , the voltage is directly driven by white noise, and  $\Lambda_v$  can be expressed in powers of  $S$  and  $S^T$ . Because  $S$  satisfies its own polynomial equation (Cayley Hamilton theorem), there will be 4 coefficients for the expansion of  $S$  and 4 for  $S^T$ , resulting in 16 coefficients that define  $\Lambda_v$  for a

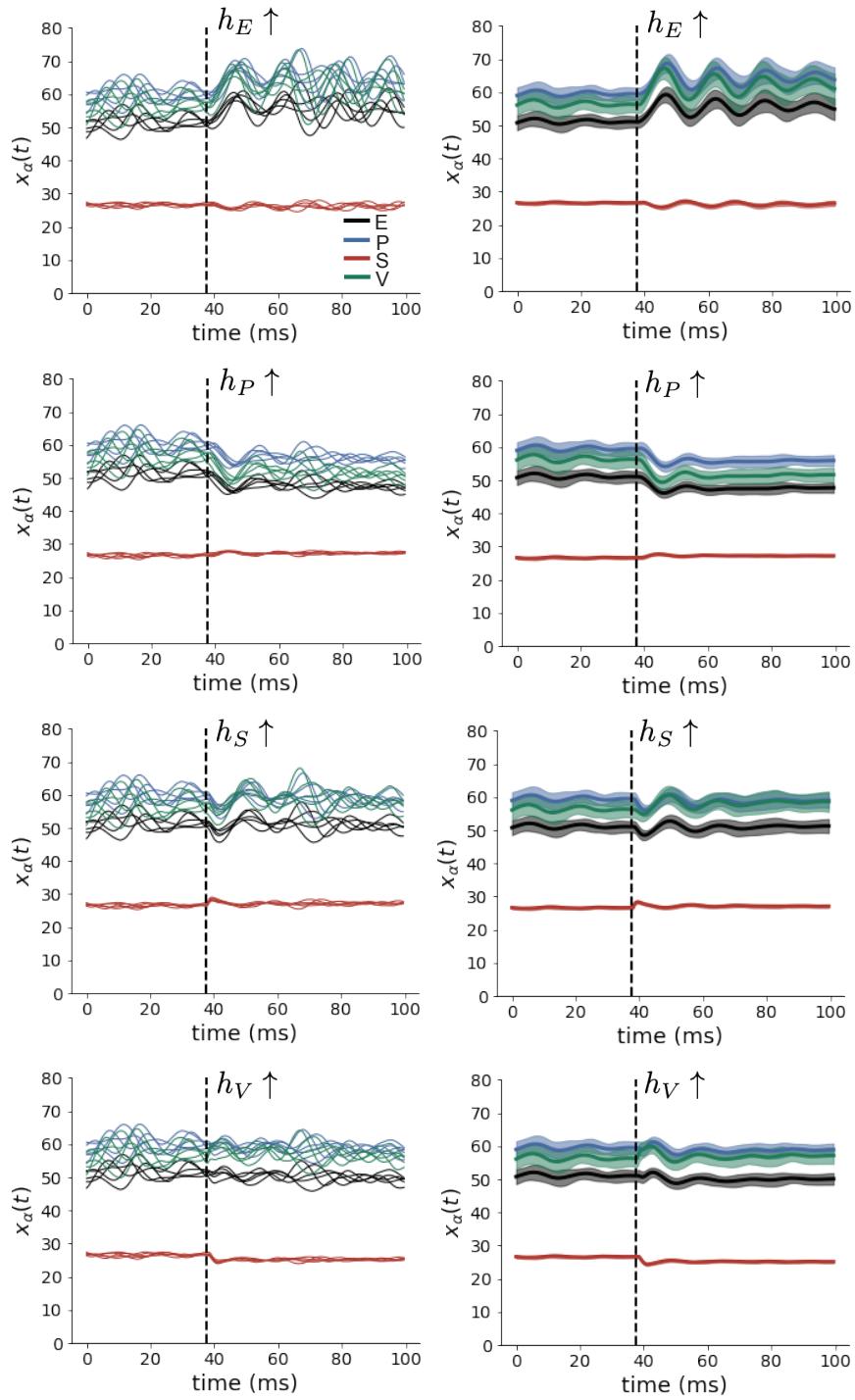


Figure S11: (Left) SSSN simulations for small increases in neuron-type population input. (Right) Average (solid) and standard deviation (shaded) of stochastic fluctuations of responses.

given  $S$ . Due to symmetry arguments [66], in this case the diagonal elements of the covariance matrix of the voltage will have the form:

$$\Lambda_{v_{ii}} = \sum_{i=\{E,P,S,V\}} g_i(S) \sigma_{ii}^2 \quad (84)$$

1282 These coefficients  $g_i(S)$  are complicated functions of the Jacobian of the system. Although expres-  
 1283 sions for these coefficients can be found explicitly, only numerical evaluation of those expressions  
 1284 determine which components of the noisy input are going to strongly influence the variability of ex-  
 1285 citatory population. Showing the generality of this dependence in more complicated noise scenarios  
 1286 (e.g.  $\tau_{\text{noise}} > \tau$  as in Section 3.4), is the focus of current research.

## 1287 5.5 Superior colliculus

### 1288 5.5.1 SC model

1289 The ability to switch between two separate tasks throughout randomly interleaved trials, or “rapid  
 1290 task switching,” has been studied in rats, and midbrain superior colliculus (SC) has been show to  
 1291 play an important in this computation [67]. Neural recordings in SC exhibited two populations of  
 1292 neurons that simultaneously represented both task context (Pro or Anti) and motor response (con-  
 1293 tralateral or ipsilateral to the recorded side), which led to the distinction of two functional classes:  
 1294 the Pro/Contra and Anti/Ipsi neurons [48]. Given this evidence, Duan et al. proposed a model  
 1295 with four functionally-defined neuron-type populations: two in each hemisphere corresponding to  
 1296 the Pro/Contra and Anti/Ipsi populations. We study how the connectivity of this neural circuit  
 1297 governs rapid task switching ability.

1298 The four populations of this model are denoted as left Pro (LP), left Anti (LA), right Pro (RP)  
 1299 and right Anti (RA). Each unit has an activity ( $x_\alpha$ ) and internal variable ( $u_\alpha$ ) related by

$$x_\alpha = \phi(u_\alpha) = \left( \frac{1}{2} \tanh \left( \frac{u_\alpha - a}{b} \right) + \frac{1}{2} \right), \quad (85)$$

1300 where  $\alpha \in \{LP, LA, RA, RP\}$ ,  $a = 0.05$  and  $b = 0.5$  control the position and shape of the nonlin-  
 1301 earity. We order the neural populations of  $x$  and  $u$  in the following manner

$$\mathbf{x} = \begin{bmatrix} x_{LP} \\ x_{LA} \\ x_{RP} \\ x_{RA} \end{bmatrix} \quad \mathbf{u} = \begin{bmatrix} u_{LP} \\ u_{LA} \\ u_{RP} \\ u_{RA} \end{bmatrix}, \quad (86)$$

1302 which evolve according to

$$\tau \frac{d\mathbf{u}}{dt} = -\mathbf{u} + W\mathbf{x} + \mathbf{h} + d\mathbf{B}. \quad (87)$$

1303 with time constant  $\tau = 0.09s$ , step size 24ms and Gaussian noise  $d\mathbf{B}$  of variance  $0.2^2$ . These  
1304 hyperparameter values are motivated by modeling choices and results from [48].

1305 The weight matrix has 4 parameters for self  $sW$ , vertical  $vW$ , horizontal  $hW$ , and diagonal  $dW$   
1306 connections:

$$W = \begin{bmatrix} sW & vW & hW & dW \\ vW & sW & dW & hW \\ hW & dW & sW & vW \\ dW & hW & vW & sW \end{bmatrix}. \quad (88)$$

1307 We study the role of parameters  $\mathbf{z} = [sW, vW, hW, dW]^\top$  in rapid task switching.

1308 The circuit receives four different inputs throughout each trial, which has a total length of 1.8s.

$$\mathbf{h} = \mathbf{h}_{\text{constant}} + \mathbf{h}_{\text{P,bias}} + \mathbf{h}_{\text{rule}} + \mathbf{h}_{\text{choice-period}} + \mathbf{h}_{\text{light}}. \quad (89)$$

1309 There is a constant input to every population,

$$\mathbf{h}_{\text{constant}} = I_{\text{constant}}[1, 1, 1, 1]^\top, \quad (90)$$

1310 a bias to the Pro populations

$$\mathbf{h}_{\text{P,bias}} = I_{\text{P,bias}}[1, 0, 1, 0]^\top, \quad (91)$$

1311 rule-based input depending on the condition

$$\mathbf{h}_{\text{P,rule}}(t) = \begin{cases} I_{\text{P,rule}}[1, 0, 1, 0]^\top, & \text{if } t \leq 1.2s \\ 0, & \text{otherwise} \end{cases} \quad (92)$$

1312

$$\mathbf{h}_{\text{A,rule}}(t) = \begin{cases} I_{\text{A,rule}}[0, 1, 0, 1]^\top, & \text{if } t \leq 1.2s \\ 0, & \text{otherwise} \end{cases}, \quad (93)$$

1313 a choice-period input

$$\mathbf{h}_{\text{choice}}(t) = \begin{cases} I_{\text{choice}}[1, 1, 1, 1]^\top, & \text{if } t > 1.2s \\ 0, & \text{otherwise} \end{cases}, \quad (94)$$

1314 and an input to the right or left-side depending on where the light stimulus is delivered

$$\mathbf{h}_{\text{light}}(t) = \begin{cases} I_{\text{light}}[1, 1, 0, 0]^\top, & \text{if } 1.2s < t < 1.5s \text{ and Left} \\ I_{\text{light}}[0, 0, 1, 1]^\top, & \text{if } 1.2s < t < 1.5s \text{ and Right} \\ 0, & \text{otherwise} \end{cases}. \quad (95)$$

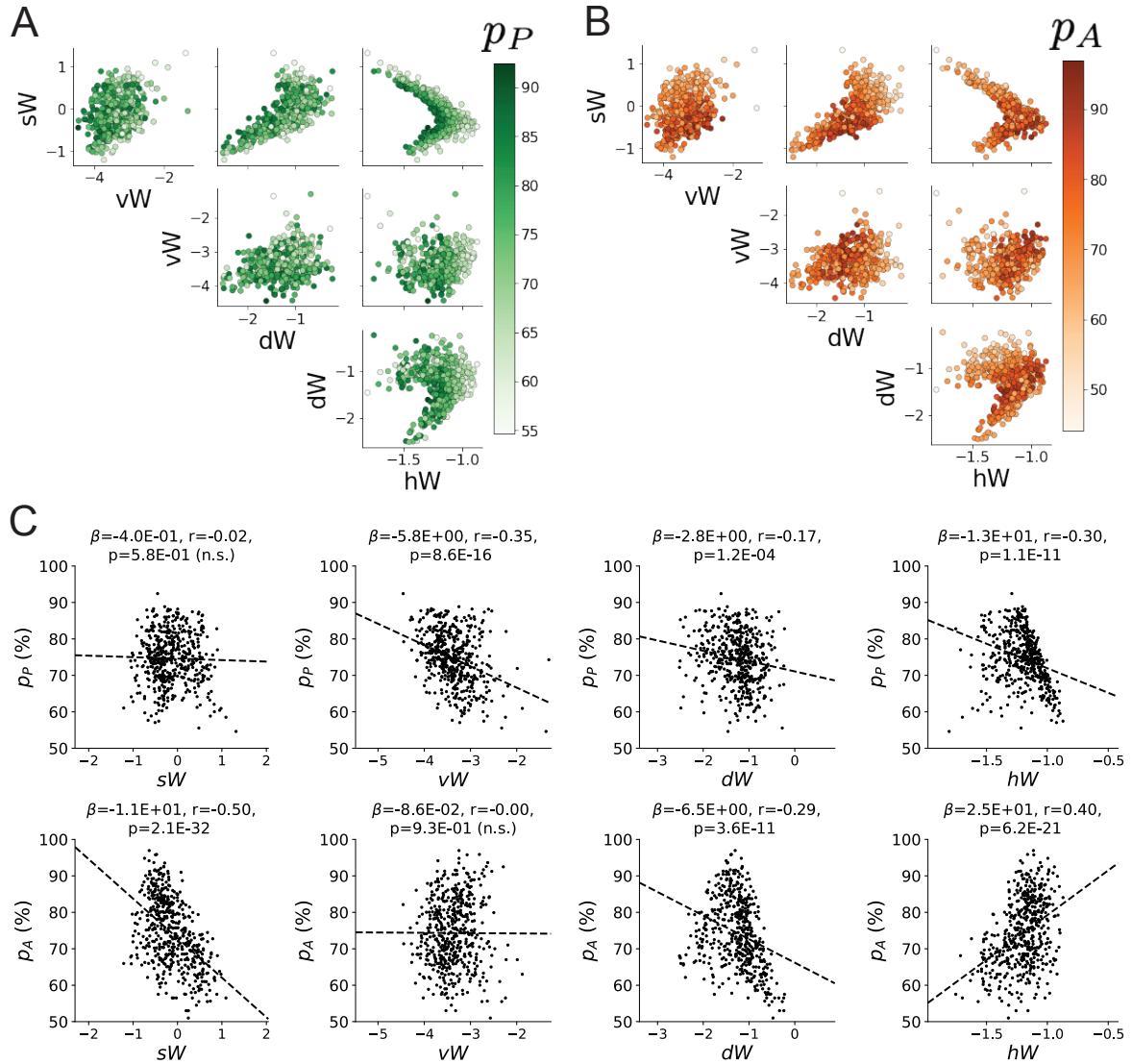


Figure S12: **A.** Same pairplot as Fig. 4C colored by Pro task accuracy. **B.** Same as A colored by Anti task accuracy. **C.** Connectivity parameters of EPI distributions versus task accuracies.  $\beta$  is slope coefficient of linear regression,  $r$  is correlation, and  $p$  is the two-tailed p-value.

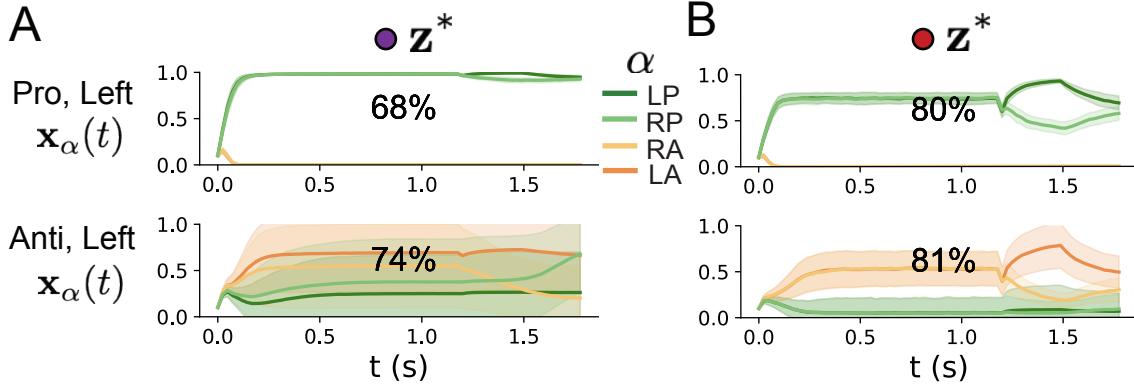


Figure S13: **A.** Simulations in network regime 1 ( $hW_{\text{fixed}} = -1.5$ ). **B.** Simulations in network regime 2 ( $hW_{\text{fixed}} = -1.5$ ).

<sup>1315</sup> The input parameterization was fixed to  $I_{\text{constant}} = 0.75$ ,  $I_{P,\text{bias}} = 0.5$ ,  $I_{P,\text{rule}} = 0.6$ ,  $I_{A,\text{rule}} = 0.6$ ,  
<sup>1316</sup>  $I_{\text{choice}} = 0.25$ , and  $I_{\text{light}} = 0.5$ .

### <sup>1317</sup> 5.5.2 Task accuracy calculation

<sup>1318</sup> The accuracies of each Pro and Anti tasks are calculated as

$$p_P(\mathbf{x}; \mathbf{z}) = \mathbb{E}_{\mathbf{x}} [\Theta[x_{LP}(t = 1.8s) - x_{RP}(t = 1.8s)]] \quad (96)$$

<sup>1319</sup> and

$$p_A(\mathbf{x}; \mathbf{z}) = \mathbb{E}_{\mathbf{x}} [\Theta[x_{RP}(t = 1.8s) - x_{LP}(t = 1.8s)]] \quad (97)$$

<sup>1320</sup> given that the stimulus is on the left side and  $\Theta$  approximates the Heaviside step function. Our  
<sup>1321</sup> accuracy calculation only considers one stimulus presentation (Left), since the model is left-right  
<sup>1322</sup> symmetric. The accuracy is averaged over 200 independent trials, and the Heaviside step function  
<sup>1323</sup> is approximated as

$$\Theta(\mathbf{x}) = \text{sigmoid}(\beta_{\Theta}\mathbf{x}), \quad (98)$$

<sup>1324</sup> where  $\beta_{\Theta} = 100$ .

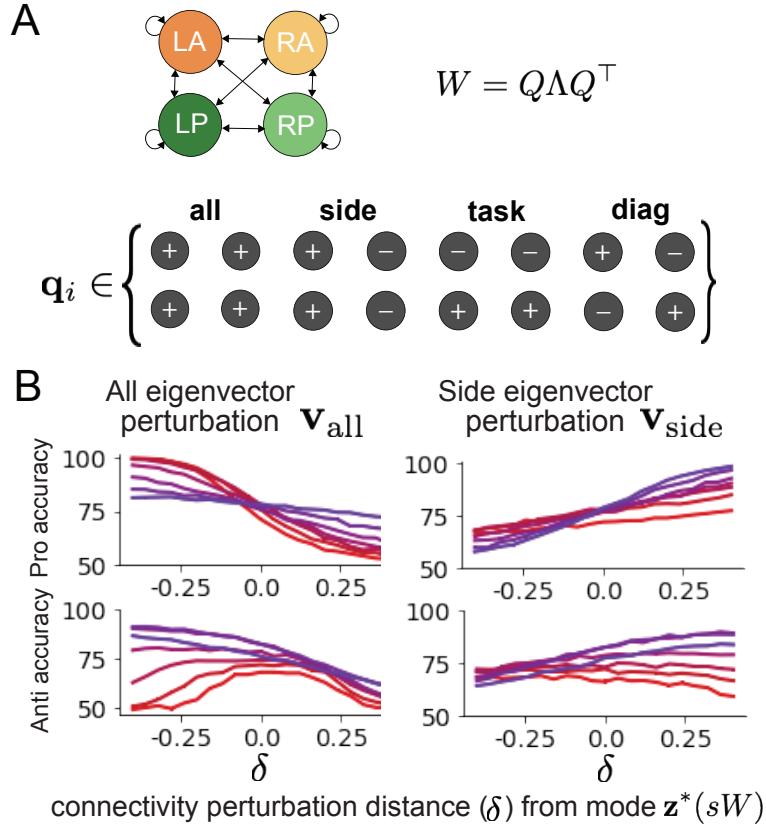


Figure S14: **A.** Invariant eigenvectors of connectivity matrix  $W$ . **B.** Accuracies for connectivity perturbations when changing  $\lambda_{\text{all}}$  and  $\lambda_{\text{side}}$  ( $\lambda_{\text{task}}$  and  $\lambda_{\text{diag}}$  shown in Fig. 4D).

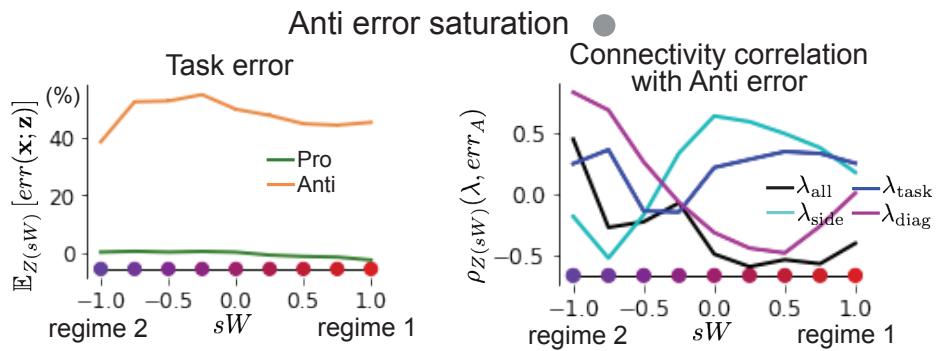


Figure S15: (Left) Mean and standard error of Pro and Anti error from regime 1 to regime 2 at  $\gamma = 0.85$ . (Right) Correlations of connectivity eigenvalues with Anti error from regime 1 to regime 2 at  $\gamma = 0.85$ .

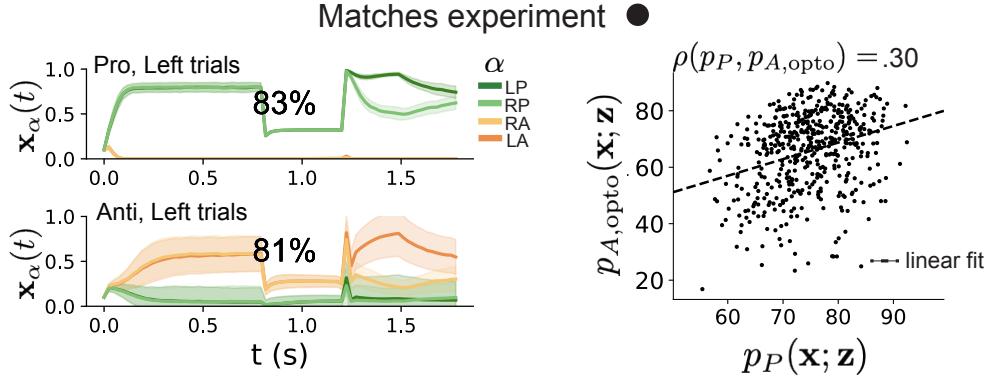


Figure S16: (Left) Mean and standard deviation (shading) of responses of the SC model at the mode of the EPI distribution to delay period inactivation at  $\gamma = 0.675$ . Accuracy in Pro (top) and Anti (bottom) task is shown as a percentage. (Right) Anti accuracy following delay period inactivation at  $\gamma = 0.675$  versus accuracy in the Pro task across connectivities in the EPI distribution.

### 1325 5.5.3 EPI details for the SC model

1326 Writing the EPI distribution as a maximum entropy distribution,  $T(\mathbf{x}, \mathbf{z})$  is comprised of both these  
1327 first and second moments of the accuracy in each task (as in Equations 20 and 21)

$$T(\mathbf{x}; \mathbf{z}) = \begin{bmatrix} p_P(\mathbf{x}; \mathbf{z}) \\ p_A(\mathbf{x}; \mathbf{z}) \\ (p_P(\mathbf{x}; \mathbf{z}) - .75)^2 \\ (p_A(\mathbf{x}; \mathbf{z}) - .75)^2 \end{bmatrix}, \quad (99)$$

$$\boldsymbol{\mu}_{\text{opt}} = \begin{bmatrix} .75 \\ .75 \\ .075^2 \\ .075^2 \end{bmatrix}. \quad (100)$$

1329 Throughout optimization, the augmented lagrangian parameters  $\eta$  and  $c$ , were updated after each  
1330 epoch of 2,000 iterations (see Section 5.1.4). The optimization converged after ten epochs (Fig.  
1331 S16).

1332 For EPI in Fig. 4C, we used a real NVP architecture with three coupling layers of affine transfor-  
1333 mations parameterized by two-layer neural networks of 50 units per layer. The initial distribution  
1334 was a standard isotropic gaussian  $\mathbf{z}_0 \sim \mathcal{N}(\mathbf{0}, I)$  mapped to a support of  $\mathbf{z}_i \in [-5, 5]$ . We used an  
1335 augmented lagrangian coefficient of  $c_0 = 10^2$ , a batch size  $n = 100$ , and  $\beta = 2$ . The distribution  
1336 was the greatest EPI distribution to converge across 5 random seeds with criteria  $N_{\text{test}} = 25$ .

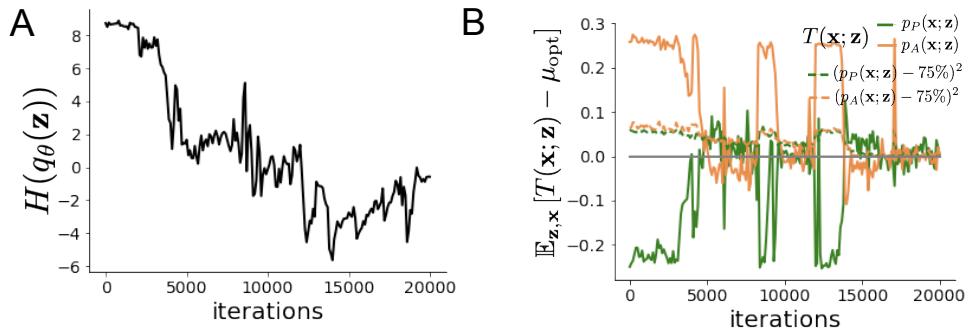


Figure S17: **A.** Entropy throughout optimization. **B.** The emergent property statistic means and variances converge to their constraints at 20,000 iterations following the tenth augmented lagrangian epoch.

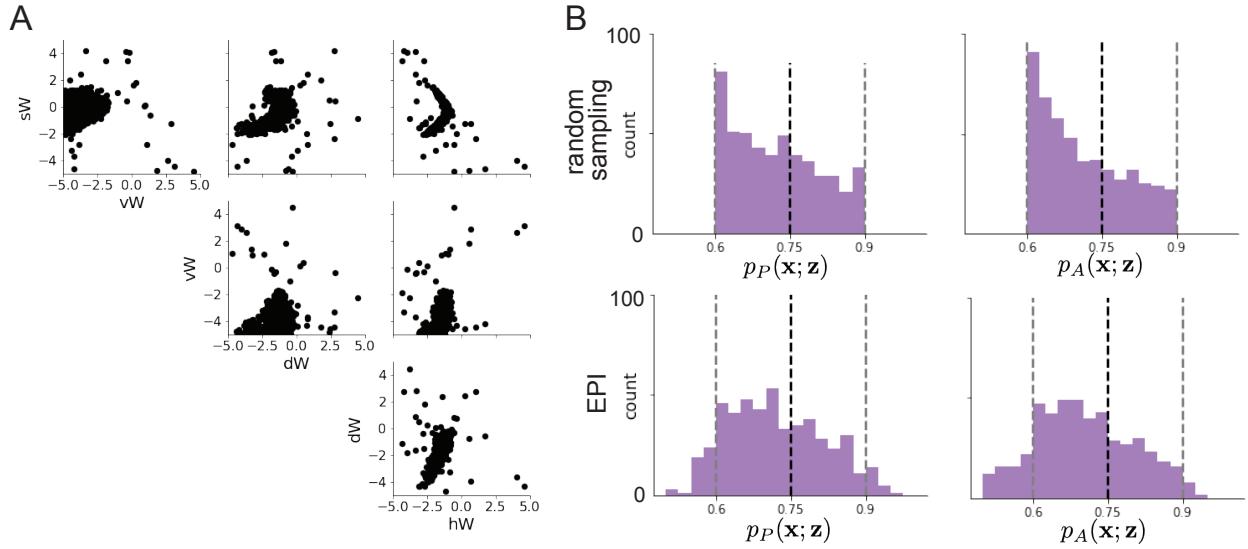


Figure S18: **A.** Entropy throughout optimization. **B.** The emergent property statistic means and variances converge to their constraints at 20,000 iterations following the tenth augmented lagrangian epoch.

1337 The bend in the EPI distribution is not a spurious result of the EPI optimization. The structure  
1338 discovered by EPI matches the shape of the set of points returned from brute-force random sampling  
1339 (Fig. S18A) These connectivities were sampled from a uniform distribution over the range of each  
1340 connectivity parameter, and all parameters producing accuracy in each task within the range of  
1341 60% to 90% were kept. This set of connectivities will not match the distribution of EPI exactly,  
1342 since it is not conditioned on the emergent property. For example the parameter set returned by  
1343 the brute-force search is biased towards lower accuracies (Fig. S18B).

#### 1344 5.5.4 Mode identification with EPI

1345 We found one mode of the EPI distribution for fixed values of  $sW$  from 1 to -1 in steps of 0.25.  
1346 To begin, we chose an initial parameter value from 500 parameter samples  $\mathbf{z} \sim q_{\theta}(\mathbf{z} \mid \mathcal{X})$  that  
1347 had closest  $sW$  value to 1. We then optimized this estimate of the mode (for fixed  $sW$ ) using  
1348 probability gradients of the deep probability distribution for 500 steps of gradient ascent with a  
1349 learning rate of  $5 \times 10^{-3}$ . The next mode (at  $sW = 0.75$ ) was found using the previous mode as  
1350 the initialization. This and all subsequent optimizations used 200 steps of gradient ascent with a  
1351 learning rate of  $1 \times 10^{-3}$ , except at  $sW = -1$  where a learning rate of  $5 \times 10^{-4}$  was used. During all  
1352 mode identification optimizations, the learning rate was reduced by half (decay = 0.5) after every  
1353 100 iterations.

#### 1354 5.5.5 Sample grouping by mode

1355 For the analyses in Figure 5C and Figure S15, we obtained parameters for each step along the  
1356 continuum between regimes 1 and 2 by sampling from the EPI distribution. Each sample was  
1357 assigned to the closest mode  $\mathbf{z}^*(sW)$  (Equation 12). Sampling continued until 500 samples were  
1358 assigned to each mode, which took 2.67 seconds (5.34ms/sample-per-mode). It took 9.59 minutes  
1359 to obtain just 5 samples for each mode with brute force sampling requiring accuracies between 60%  
1360 and 90% in each task (115s/sample-per-mode). This corresponds to a sampling speed increase of  
1361 roughly 21,500 once the EPI distribution has been learned.

#### 1362 5.5.6 Sensitivity analysis

1363 At each mode, we measure the sensitivity dimension (that of most negative eigenvalue in the Hessian  
1364 of the EPI distribution)  $\mathbf{v}_1(\mathbf{z}^*)$ . To resolve sign degeneracy in eigenvectors, we chose  $\mathbf{v}_1(\mathbf{z}^*)$  to have

1365 negative element in  $hW$ . This tells us what parameter combination rapid task switching is most  
 1366 sensitive to at this parameter choice in the regime.

1367 **5.5.7 Connectivity eigendecomposition and processing modes**

1368 To understand the connectivity mechanisms governing task accuracy, we took the eigendecomposi-  
 1369 tion of the connectivity matrices  $W = Q\Lambda Q^{-1}$ , which results in the same eigenmodes  $\mathbf{q}_i$  for all  $W$   
 1370 parameterized by  $\mathbf{z}$  (Fig. S14A). These eigenvectors are always the same, because the connectivity  
 1371 matrix is symmetric and the model also assumes symmetry across hemispheres, but the eigenvalues  
 1372 of connectivity (or degree of eigenmode amplification) change with  $\mathbf{z}$ . These basis vectors have in-  
 1373 tuitive roles in processing for this task, and are accordingly named the *all* eigenmode - all neurons  
 1374 co-fluctuate, *side* eigenmode - one side dominates the other, *task* eigenmode - the Pro or Anti pop-  
 1375 ulations dominate the other, and *diag* mode - Pro- and Anti-populations of opposite hemispheres  
 1376 dominate the opposite pair. Due to the parametric structure of the connectivity matrix, the pa-  
 1377 rameters  $\mathbf{z}$  are a linear function of the eigenvalues  $\boldsymbol{\lambda} = [\lambda_{\text{all}}, \lambda_{\text{side}}, \lambda_{\text{task}}, \lambda_{\text{diag}}]^T$  associated with these  
 1378 eigenmodes.

$$\mathbf{z} = A\boldsymbol{\lambda} \quad (101)$$

$$1379 \quad A = \frac{1}{4} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & -1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & 1 & -1 \end{bmatrix}. \quad (102)$$

1380 We are interested in the effect of raising or lowering the amplification of each eigenmode in the  
 1381 connectivity matrix by perturbing individual eigenvalues  $\lambda$ . To test this, we calculate the unit  
 1382 vector of changes in the connectivity  $\mathbf{z}$  that result from a change in the associated eigenvalues

$$\mathbf{v}_a = \frac{\frac{\partial \mathbf{z}}{\partial \lambda_a}}{\left| \frac{\partial \mathbf{z}}{\partial \lambda_a} \right|_2}, \quad (103)$$

1383 where

$$\frac{\partial \mathbf{z}}{\partial \lambda_a} = A\mathbf{e}_a, \quad (104)$$

1384 and e.g.  $\mathbf{e}_{\text{all}} = [1, 0, 0, 0]^T$ . So  $\mathbf{v}_a$  is the normalized column of  $A$  corresponding to eigenmode  
 1385  $a$ . The parameter dimension  $\mathbf{v}_a$  ( $a \in \{\text{all}, \text{side}, \text{task}, \text{and diag}\}$ ) that increases the eigenvalue of  
 1386 connectivity  $\lambda_a$  is  $\mathbf{z}$ -invariant (Equation 104) and  $\mathbf{v}_a \perp \mathbf{v}_{b \neq a}$ . By perturbing  $\mathbf{z}$  along  $\mathbf{v}_a$ , we  
 1387 can examine how model function changes by directly modulating the connectivity amplification of  
 1388 specific eigenmodes, which having interpretable roles in processing in each task.

1389 **5.5.8 Modeling optogenetic silencing.**

1390 We tested whether the inferred SC model connectivities could reproduce experimental effects of  
1391 optogenetic inactivation in rats [48]. During periods of simulated optogenetic inactivation, activity  
1392 was decreased proportional to the optogenetic strength  $\gamma \in [0, 1]$

$$x_\alpha = (1 - \gamma)\phi(u_\alpha). \quad (105)$$

1393 Delay period inactivation was from  $0.8 < t < 1.2$ .