

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

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

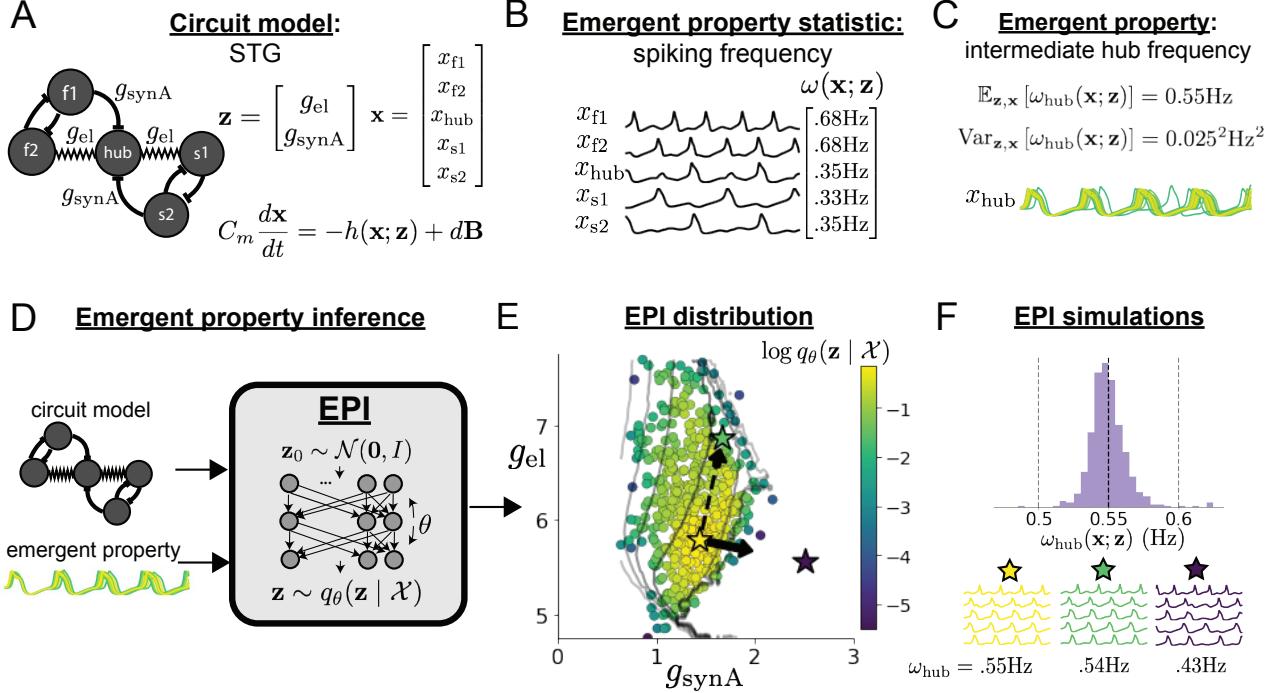


Figure 1: Emergent property inference (EPI) in the stomatogastric ganglion. **A.** Conductance-based subcircuit model of the STG. **B.** Spiking frequency  $\omega(\mathbf{x}; \mathbf{z})$  is an emergent property statistic. Simulated at  $g_{el} = 4.5\text{nS}$  and  $g_{synA} = 3\text{nS}$ . **C.** 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, ... .575 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.

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. The modes of  $q_{\theta}(\mathbf{z} | \mathcal{X})$   
165 indicate parameter choices emblematic of the emergent property (Fig. 1E yellow star). As prob-  
166 ability in the EPI distribution decreases, the emergent property deteriorates. Perturbing  $\mathbf{z}$  along  
167 a dimension in which  $q_{\theta}(\mathbf{z} | \mathcal{X})$  change little will not disturb the emergent property, making this  
168 parameter combination *robust* with respect to the emergent property. In contrast, if  $\mathbf{z}$  is perturbed  
169 along a dimension with strongly decreasing  $q_{\theta}(\mathbf{z} | \mathcal{X})$ , that parameter combination is deemed *sen-*  
170 *sitive* [27, 30]. By querying the second order derivative (Hessian) of  $\log q_{\theta}(\mathbf{z} | \mathcal{X})$  at a mode, we

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

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

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

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

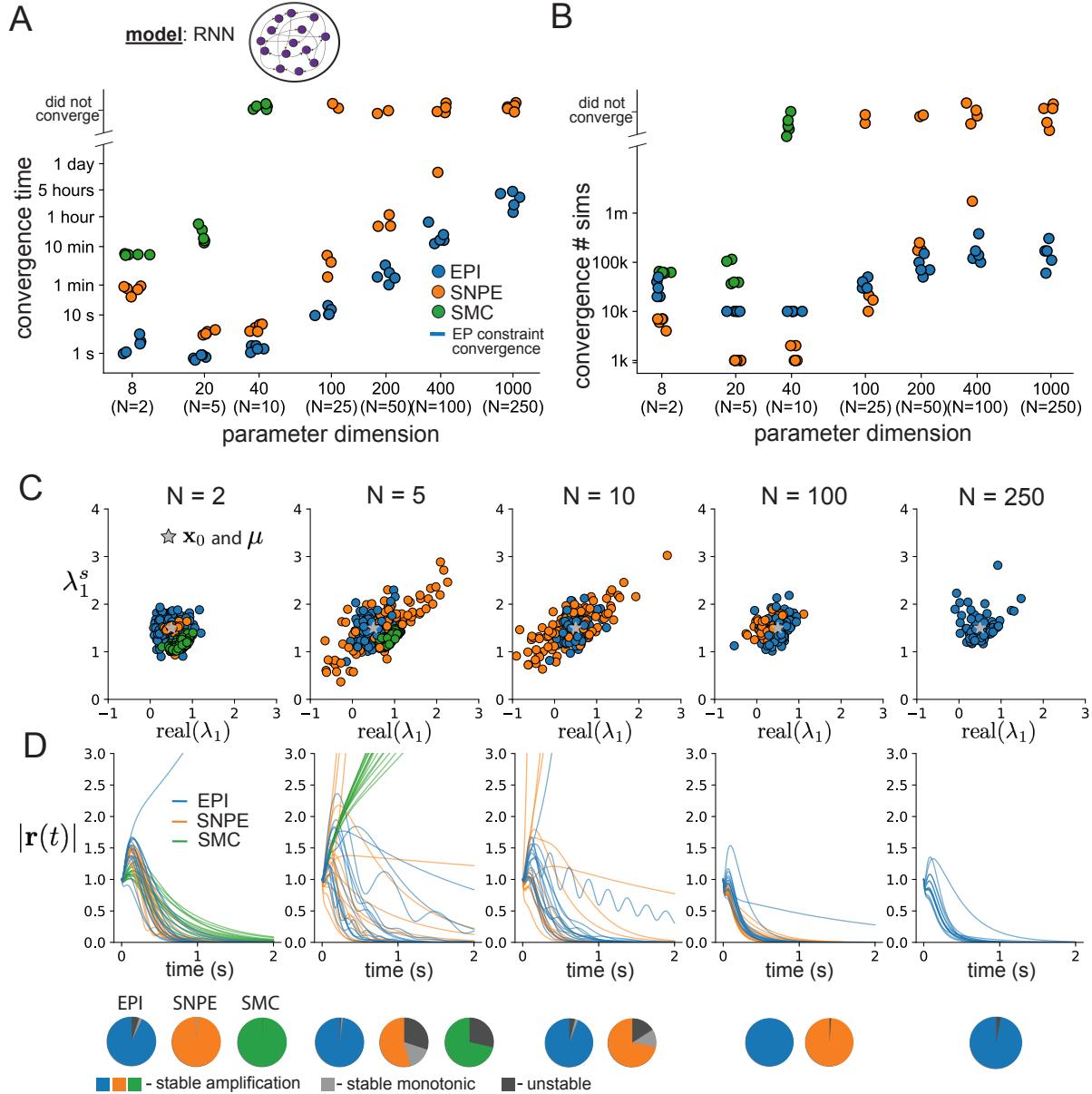


Figure 2: **A.** Wall time of EPI (blue), SNPE (orange), and SMC-ABC (green) to converge on RNN connectivities producing stable amplification. Each dot shows convergence time for an individual random seed. For reference, the mean wall time for EPI to achieve its full constraint convergence (means and 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  $\mathbf{x}_0 = \mu$  (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. 10), and tested more aggressive SNPE hyperparameterizations when SNPE failed to converge (Fig. 11). 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 and sampling efficiency (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 excitation, these models give rise to paradoxical responses [12], selective amplification [43, 50], surround

258 suppression [58] and normalization [59]. Recent theoretical work [60] shows that stabilized E-I  
 259 models reproduce the effect of variability suppression [61]. Furthermore, experimental evidence  
 260 shows that inhibition is composed of distinct elements – parvalbumin (P), somatostatin (S), VIP  
 261 (V) – composing 80% of GABAergic interneurons in V1 [63–65], and that these inhibitory cell  
 262 types follow specific connectivity patterns (Fig. 3A) [66]. Here, we use EPI on a model of V1 with  
 263 biologically realistic connectivity to show how the structure of input across neuron types affects  
 264 the variability of the excitatory population – the population largely responsible for projecting to  
 265 other brain areas [67].

266 We considered response variability of a nonlinear dynamical V1 circuit model (Fig. 3A) with a state  
 267 comprised of each neuron-type population’s rate  $\mathbf{x} = [x_E, x_P, x_S, x_V]^\top$ . Each population receives  
 268 recurrent input  $W\mathbf{x}$ , where  $W$  is the effective connectivity matrix (see Section 5.4) and an external  
 269 input with mean  $\mathbf{h}$ , which determines population rate via supralinear nonlinearity  $\phi(\cdot) = [\cdot]_+^2$ . The  
 270 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  
 271 has a slower dynamical timescale  $\tau_{\text{noise}} > \tau$  than the population rate, allowing fluctuations around  
 272 a stimulus-dependent steady-state (Fig. 3B). This model is the stochastic stabilized supralinear  
 273 network (SSSN) [60]

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

274 generalized to have multiple inhibitory neuron types. It introduces stochasticity to four neuron-  
 275 type models of V1 [46]. Stochasticity and inhibitory multiplicity introduce substantial complexity  
 276 to the mathematical treatment of this problem (see Section 5.4.4) motivating the analysis of this  
 277 model with EPI. Here, we consider fixed weights  $W$  and input  $\mathbf{h}$  [47], and study the effect of input  
 278 variability  $\mathbf{z} = [\sigma_E, \sigma_P, \sigma_S, \sigma_V]^\top$  on excitatory variability.

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

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

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

283 First, we ran EPI to obtain parameter distribution  $q_{\theta}(\mathbf{z} | \mathcal{X}(5 \text{ Hz}))$  producing E-population vari-  
 284 ability around 5 Hz (Fig. 3D). From the marginal distribution of  $\sigma_E$  and  $\sigma_P$  (Fig. 3D, top-left),  
 285 we can see that  $s_E(\mathbf{x}; \mathbf{z})$  is sensitive to various combinations of  $\sigma_E$  and  $\sigma_P$ . Alternatively, both  $\sigma_S$

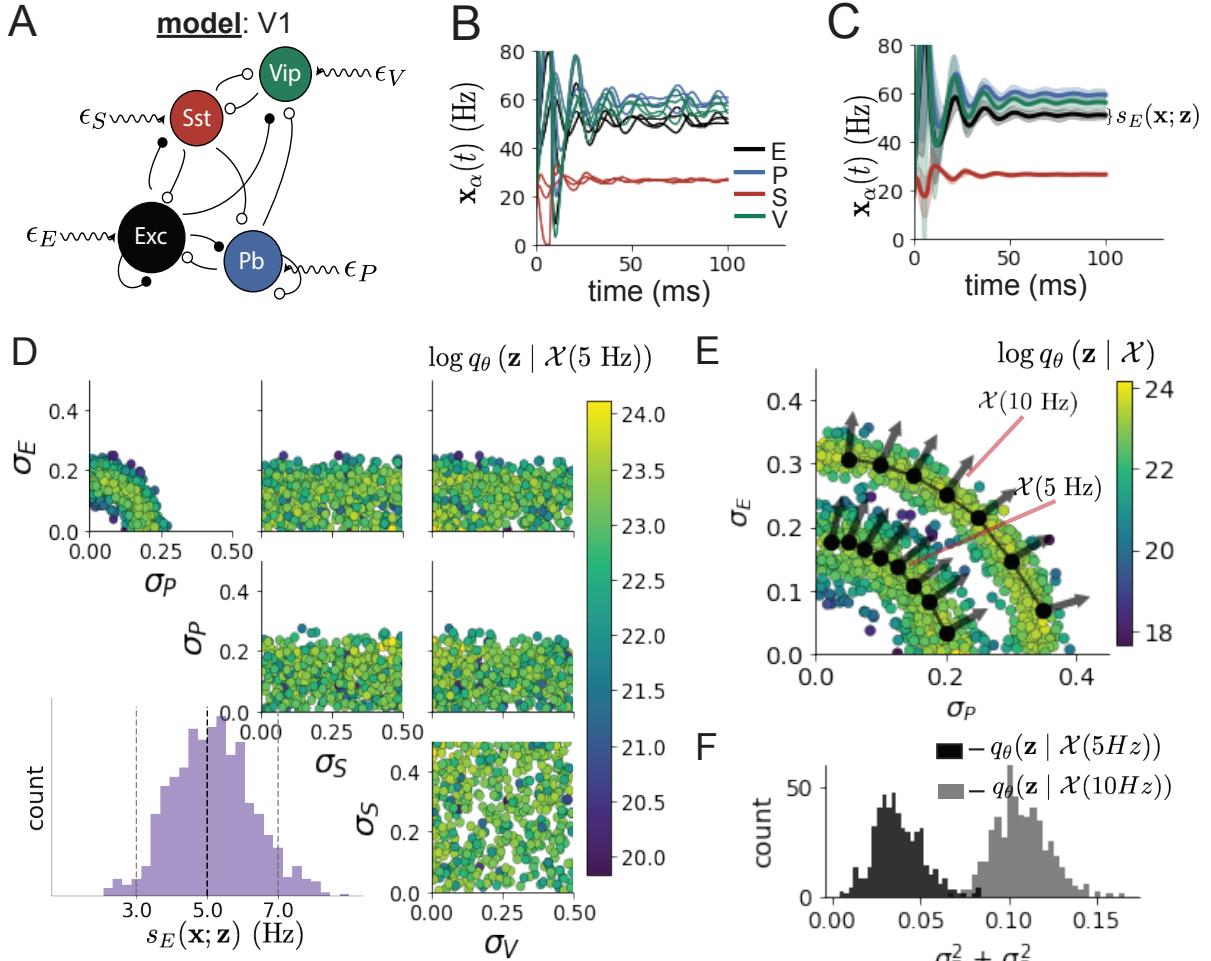


Figure 3: Emergent property inference in the stochastic stabilized supralinear network (SSSN)

**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 shown 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 and  $\sigma_V$  are degenerate with respect to  $s_E(\mathbf{x}; \mathbf{z})$  evidenced by the, unexpectedly, high variability in  
287 those dimensions (Fig. 3D, bottom-right). Together, these observations imply a curved path with  
288 respect to  $s_E(\mathbf{x}; \mathbf{z})$  of 5 Hz, which is indicated by the modes along  $\sigma_P$  (Fig. 3E).

289 Figure 3E suggests a quadratic relationship in E-population fluctuations and the standard deviation  
290 of E- and P-population input; as the square of either  $\sigma_E$  or  $\sigma_P$  increases, the other compensates  
291 by decreasing to preserve the level of  $s_E(\mathbf{x}; \mathbf{z})$ . This quadratic relationship is preserved at greater  
292 level of E-population variability  $\mathcal{X}(10 \text{ Hz})$  (Fig. 3E). Indeed, the sum of squares of  $\sigma_E$  and  $\sigma_P$  is  
293 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  
294 of  $\sigma_S$  and  $\sigma_V$  are not significantly different in the two EPI distributions (Fig. 15,  $p = .40$ ), in which  
295 parameters were bounded from 0 to 0.5. The strong interaction between E- and P-population input  
296 variability on excitatory variability is intriguing, since this circuit exhibits a paradoxical effect in  
297 the P-population (and no other inhibitory types) (Fig. 15), meaning that the E-population is P-  
298 stabilized. Future research may uncover a link between the population of network stabilization and  
299 compensatory interactions governing excitatory variability.

300 EPI revealed the quadratic dependence of excitatory variability on input variability to the E- and  
301 P-populations, as well as its independence to input from the other two inhibitory populations. In a  
302 simplified model ( $\tau = \tau_{\text{noise}}$ ), it can be shown that surfaces of equal variance are ellipsoids as a func-  
303 tion of  $\sigma$  (see Section 5.4.4). Nevertheless, the sensitive and degenerate parameters are challenging  
304 to predict mathematically, since the covariance matrix depends on the steady-state solution of the  
305 network [60, 68], and terms in the covariance expression increase quadratically with each additional  
306 neuron-type population (see also Section 5.4.4). By pointing out this mathematical complexity, we  
307 emphasize the value of streamlined methods for gaining understanding about theoretical models  
308 when mathematical analysis becomes onerous or impractical. While we have just shown that EPI  
309 can be used to investigate fundamental aspects of sensory computation, in the next two sections,  
310 we use the probabilistic tools of EPI to identify and characterize two distinct parametric regimes of  
311 a neural circuit executing a computation, and then relate these insights to behavioral experiments.

### 312 3.5 EPI identifies two regimes of rapid task switching

313 It has been shown that rats can learn to switch from one behavioral task to the next on randomly  
314 interleaved trials [69], and an important question is what types of neural connectivity produce this  
315 computation. In this experimental setup, rats were explicitly cued on each trial to either orient  
316 towards a visual stimulus in the Pro (P) task or orient away from the stimulus in the Anti (A)

task (Fig. 4A). Neural recordings in superior colliculus (SC) exhibited two populations of neurons that represented task context (Pro or Anti). Furthermore, Pro/Anti neurons in each hemisphere were strongly correlated with the animal’s decision [48]. These results motivated a model of SC that is a four-population dynamical system with functionally-defined neuron-types. Here, our goal is to understand how connectivity in this circuit model governs the ability to perform rapid task switching: to respond with satisfactory accuracy in both tasks on randomly interleaved trials.

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

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

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

Rapid task switching is formalized mathematically as an emergent property with two statistics: 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 .75 in each task with variance  $.075^2$

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

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

The EPI inferred distribution (Fig. 4C) produces Pro and Anti task accuracies (Fig. 4C, middle-left) consistent with rapid task switching (Equation 9). This parameter distribution has rich structure, that is not captured well by simple linear correlations (Fig. 17). Specifically, the shape of the EPI distribution is sharply bent, matching ground truth structure indicated by brute-force

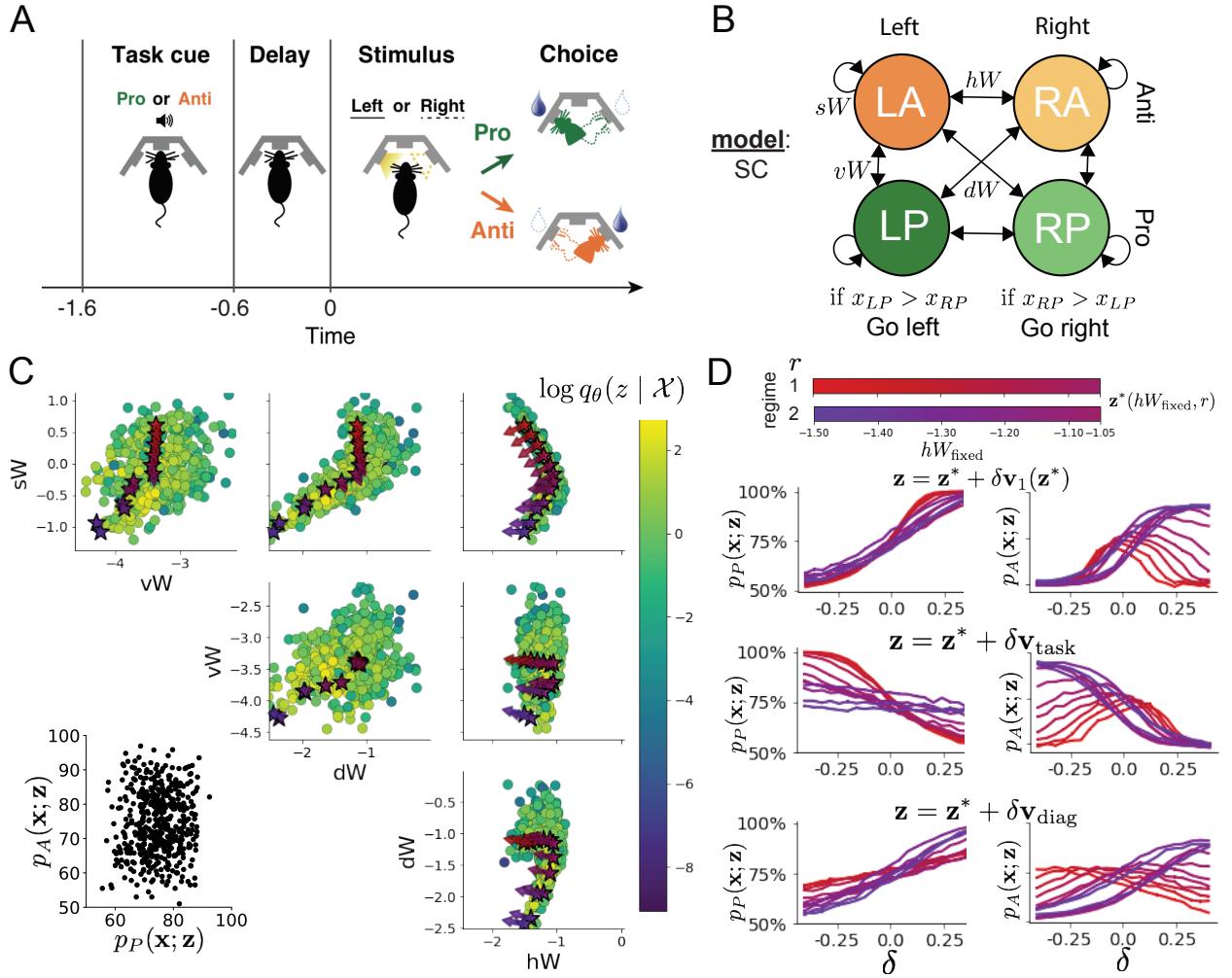


Figure 4: **A.** Rapid task switching behavioral paradigm (see text). **B.** Model of superior colliculus (SC). Neurons: LP - Left Pro, RP - Right Pro, LA - Left Anti, RA - Right Anti. Parameters:  $sW$  - self,  $hW$  - horizontal,  $vW$  - vertical,  $dW$  - diagonal weights. **C.** The EPI inferred distribution of rapid task switching networks. Stars indicate modes  $\mathbf{z}^*$  whose color indicates connectivity regime (see legend Fig 4D). Sensitivity vectors  $\mathbf{v}_1(\mathbf{z}^*)$  are shown by arrows. (Bottom-left) EPI predictive distribution of task accuracies. **D.** The connectivity regimes have different responses to perturbation. (Top) Mean and standard error ( $N_{\text{test}} = 25$ ) of accuracy with respect to perturbation along the sensitivity dimension of each mode  $\mathbf{z}^*$ . (Middle) Same with perturbation in the dimension of increasing  $\lambda_{\text{task}}$  ( $\mathbf{v}_{\text{task}}$ ). (Bottom) Same with perturbation in the dimension of increasing  $\lambda_{\text{diag}}$  ( $\mathbf{v}_{\text{diag}}$ ).

342 sampling (Fig. 24) This is most saliently pointed out in the marginal distribution of  $sW-hW$  (Fig.  
343 4C top-right), where anticorrelation between  $sW$  and  $hW$  switches to correlation with decreasing  
344  $sW$ . The two regimes produce different types of responses in the Pro and Anti tasks (Fig. SC2).  
345 Not only has EPI captured this complicated distribution of connectivities producing rapid task  
346 switching, we can query the EPI distribution  $q_{\theta}(\mathbf{z} | \mathcal{X})$  to understand these two parametric regimes  
347 of SC connectivity.

348 To distinguish these two regimes, we use the EPI distribution to identify two sets of modes. By  
349 fixing  $hW$  to different values and doing gradient ascent on  $\log q_{\theta}(\mathbf{z} | \mathcal{X})$ , we arrive at two solutions  
350  $\mathbf{z}^*(hW_{\text{fixed}}, r)$  where regime  $r \in [1, 2]$ , and regime 1 is that of greater  $sW$  (see Section 5.5.4). As  
351  $hW_{\text{fixed}}$  increases, the modes coalesce to intermediate parameters reflecting a transition between  
352 the two sets of modes (Fig. 20 top). By using EPI to connect these two regimes through this  
353 transitional region of parameter space, we can explore what distinguishes the two regimes by  
354 stepping from the prototypical connectivity of regime 1 to that of regime 2.

355 Although the connectivities gradually coalesce to the transitional part of parameter space, the  
356 sensitivity dimensions  $\mathbf{v}_1(\mathbf{z})$  are categorically different across regimes (Fig. 20 bottom). The  
357 sensitivity dimension identifies the parameter combination which causes the emergent property to  
358 diminish with the shortest perturbation. Since the two regimes have different  $\mathbf{v}_1(\mathbf{z})$ , this suggests  
359 they have different pathologies in their connectivity. By perturbing connectivity in each regime  
360 along the sensitivity dimension, we can get a sense of the differing nature of these pathologies.

361 When perturbing connectivity along the sensitivity dimension, Pro accuracy monotonically in-  
362 creases in both regimes (Fig. 4D, top-left). However, there is a stark difference between regimes in  
363 Anti accuracy. Anti accuracy falls in either direction of  $\mathbf{v}_1$  in regime 1, yet monotonically increases  
364 along with Pro accuracy in regime 2 (Fig. 4D, top-right). These distinct pathologies of rapid task  
365 switching are caused by distinct connectivity changes ( $\mathbf{v}_1(\mathbf{z}^*(r=1))$  vs  $\mathbf{v}_1(\mathbf{z}^*(r=2))$ ) and explain  
366 the sharp change in local structure of the EPI distribution.

367 To further examine the two regimes, we can perturb  $\mathbf{z}$  in the same way along dimensions that inde-  
368 pendently change the eigenvalues of the connectivity matrix (which has constant eigenvectors with  
369 respect to  $\mathbf{z}$ ). These eigenvalues  $\lambda_{\text{all}}$ ,  $\lambda_{\text{side}}$ ,  $\lambda_{\text{task}}$ , and  $\lambda_{\text{diag}}$  correspond to connectivity eigenmodes  
370 with intuitive roles in processing in this task (Fig. 19A). For example, greater  $\lambda_{\text{task}}$  will strengthen  
371 internal representations of task, while greater  $\lambda_{\text{diag}}$  will amplify dominance of Pro and Anti pairs in  
372 opposite hemispheres (Section 5.5.6). Perturbation analyses reveal that decreasing  $\lambda_{\text{task}}$  has a very  
373 similar effect on Anti accuracy as perturbations along the sensitivity dimension (Fig. 4D, middle).

374 This suggests that there is a carefully tuned strength of task representation in connectivity regime  
375 1, which if disturbed results in random Anti trial responses. Finally, we recognize that increasing  
376  $\lambda_{\text{diag}}$  has opposite effects on Anti accuracy in each regime (Fig. 4D, bottom). In the next section,  
377 we build on these mechanistic characterizations of each regime by examining their resilience to  
378 optogenetic silencing.

379 **3.6 EPI inferred SC connectivities reproduce results from optogenetic inacti-  
380 vation experiments**

381 During the delay period of this task, the circuit must prepare to execute the correct task according  
382 to the presented cue. Experimental results from Duan et al. found that bilateral optogenetic  
383 inactivation of SC during the delay period consistently decreased performance in the Anti task, but  
384 had no effect on the Pro task (Fig. 5A). This suggests that SC maintains a representation of task  
385 throughout the delay period, which is important for correct execution of the Anti task. Network  
386 connectivities inferred by EPI exhibited this same effect in simulation at high optogenetic strengths  
387  $\gamma$  (Fig. 5B) (see Section 5.5.7).

388 The mean increase in Anti error is closest to the experimentally measured value of 7% at  $\gamma = 0.675$   
389 (Fig. 5B, black dot). At this level of optogenetic strength, only regime 1 exhibits an increase in  
390 Anti error with delay period silencing (Fig. 5C, left). The connectivities in regime 2 are thus more  
391 resilient to delay period silencing during Anti trials than regime 1. In regime 1, greater  $\lambda_{\text{task}}$  and  
392  $\lambda_{\text{diag}}$  decrease Anti error (Fig. 5C, right). In other words, these anticorrelations show that stronger  
393 task representations and diagonal amplification make the SC model more resilient to delay period  
394 silencing in the Anti task. All correlations of connectivity eigenvalue with Anti error degrade in  
395 regime 2, where there is no effect of delay period silencing on Anti error (Fig. 5C, right).

396 At about  $\gamma = 0.85$  (Fig. 5B, gray dot), the Anti error saturates, while Pro error remains at zero  
397 Following delay period inactivation at this optogenetic strength, there are strong similarities in  
398 the responses of Pro and Anti trials during the choice period (Fig. 5D, left). We interpreted  
399 these similarities to suggest that delay period inactivation at this saturated level flips the internal  
400 representation of task (from Anti to Pro) in the circuit model. This would explain why the Anti  
401 error saturates at 50%: the average Anti accuracy in EPI inferred connectivities is 75%, but is 25%  
402 when the internal representation is flipped during delay period silencing. This hypothesis prescribes  
403 a model of Anti accuracy during delay period silencing of  $p_{A,\text{opto}} = 100\% - p_P$ , which is fit closely  
404 across both regimes of the EPI inferred connectivities (Fig. 5D, right). Similarities between Pro

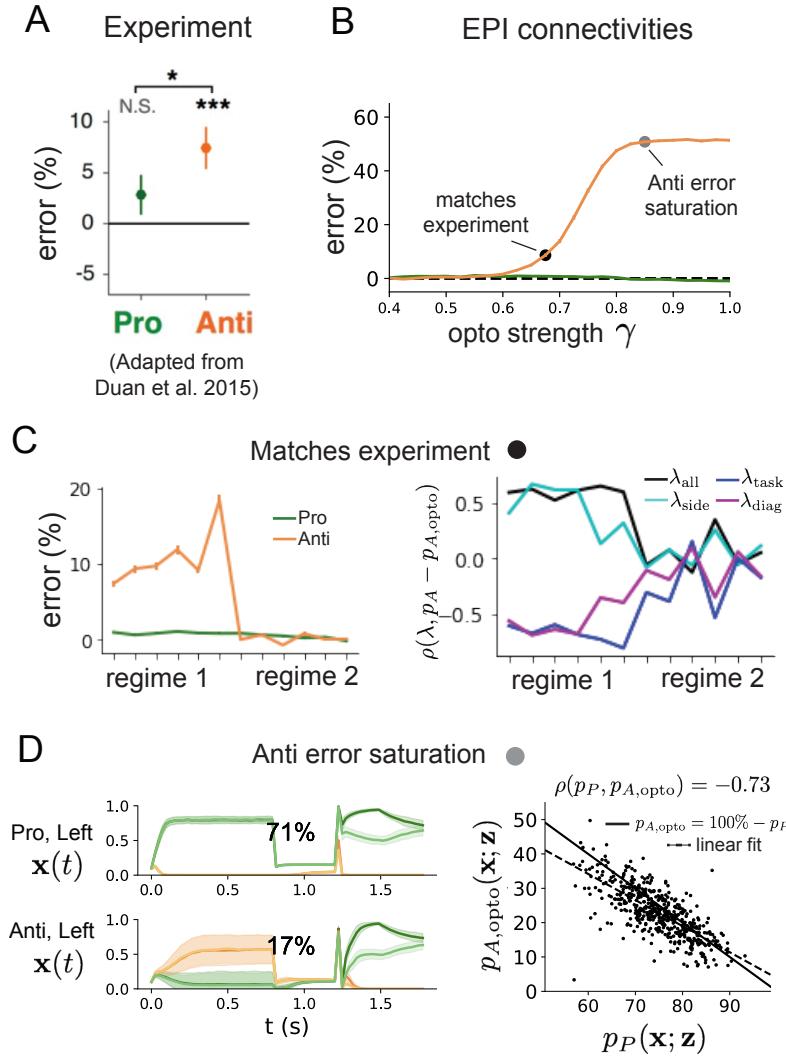


Figure 5: **A.** Experimental effect of delay period silencing on Pro and Anti task accuracy in rats. **B.** Mean and standard error (bars) of task error induced by delay period inactivation of varying optogenetic strength across the EPI distribution. **C.** (Left) Mean and standard error of Pro and 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  $\gamma = 0.85$ . (Right) Anti accuracy following delay period inactivation at  $\gamma = 0.85$  versus accuracy in the Pro task across connectivities in the EPI distribution.

405 and Anti trial responses were not present at the experiment-matching level of  $\gamma = 0.675$  (Fig. 22  
406 left) and neither was anti-correlation in  $p_P$  and  $p_{A,\text{opto}}$  (Fig. 22 right).

407 In summary, the connectivity inferred by EPI to perform rapid task switching replicated results  
408 from optogenetic silencing experiments. We found that at levels of optogenetic strength matching  
409 experimental levels of Anti error, only one regime actually exhibited the effect. This suggests that  
410 one regime is less resilient to optogenetic perturbation, and perhaps more biologically realistic.  
411 Finally, we mechanistically characterized the pathology in Anti error that occurs in both regimes  
412 when optogenetic strength is increased to high levels. The probabilistic tools afforded by EPI  
413 yielded this insight: we identified two regimes and the continuum of connectivities between them  
414 by taking gradients of parameter probabilities in the EPI distribution, we identified sensitivity  
415 dimensions by measuring the Hessian of the EPI distribution, and we obtained many parameter  
416 samples at each step along the continuum (in 7.36 seconds with the EPI distribution rather than  
417 4.2 days with brute force methods, see Section 5.5).

## 418 4 Discussion

419 In neuroscience, machine learning has primarily been used to reveal structure in neural datasets [20].  
420 Careful inference procedures are developed for these statistical models allowing precise, quantita-  
421 tive reasoning, which clarifies the way data informs beliefs about the model parameters. However,  
422 these statistical models often lack resemblance to the underlying biology, making it unclear how  
423 to go from the structure revealed by these methods, to the neural mechanisms giving rise to it. In  
424 contrast, theoretical neuroscience has primarily focused on careful models of neural circuits and  
425 the production of emergent properties of computation, rather than measuring structure in neural  
426 datasets. In this work, we improve upon parameter inference techniques in theoretical neuroscience  
427 with emergent property inference, harnessing deep learning towards parameter inference with re-  
428 spect to computation in neural circuit models (see Section 5.1.1).

429 Methodology for statistical inference in circuit models has evolved considerably in recent years.  
430 Early work used rejection sampling techniques [24–26], but EPI and other recently developed  
431 methodology (SNPE [35]) employ deep learning to improve efficiency and provide deep, flexible  
432 distribution approximations. SNPE has been used for posterior inference of parameters in circuit  
433 models conditioned upon exemplar data used to represent computation, but it does not infer param-  
434 eter distributions that only produce the computation of interest like EPI (see Section 3.3). Finally,

435 we show that EPI has better scaling properties than SNPE when emergent property gradients are  
436 tractable (Section 3.3). In summary, choice of deep inference technique should consider emergent  
437 property complexity and differentiability, dimensionality of parameter space, and the importance  
438 of constraining the model behavior predicted by the inferred parameter distribution.

439 In this paper, we prove out the value of deep inference for parameter sensitivity analyses at both the  
440 local and global level. With these techniques, flexible deep probability distributions are optimized to  
441 capture global structure by approximating the full distribution of suitable parameters. Importantly,  
442 the local structure of this deep probability distribution can be quantified at any parameter choice,  
443 offering instant sensitivity measurements after fitting. For example, the global structure captured  
444 by EPI revealed two distinct parameter regimes, which had different local structure quantified by  
445 the deep probability distribution (see Section 5.5). In comparison, Bayesian MCMC is considered a  
446 popular approach for capturing global parameter structure [70], but there is no variational approx-  
447 imation (the deep probability distribution in EPI), so sensitivity information is not queryable and  
448 sampling remains slow after convergence. Local sensitivity analyses (e.g. [27]) may be performed  
449 independently at individual parameter samples, but these methods alone do not capture the full  
450 picture in nonlinear, complex distributions. Therefore, deep inference is the only approach yielding  
451 an object – the deep probability distribution – which produces a wholistic assessment of parameter  
452 sensitivity at the local and global level, which were used in this study to make novel insights into  
453 theoretical models of neural computation. Together, the abilities to condition upon emergent prop-  
454 erties, the efficient inference algorithm, and the capacity for parameter sensitivity analyses make  
455 EPI a powerful new method for addressing inverse problems in theoretical neuroscience.

456 Even with a high degree of biophysical realism and expensive emergent property gradients, EPI was  
457 run successfully on intermediate hub frequency in a 5-neuron subcircuit model of the STG (Section  
458 3.1). However, conditioning on the pyloric rhythm [57] in a model of the pyloric subnetwork  
459 model [15] proved to be prohibitive with EPI. The pyloric subnetwork requires many time steps for  
460 simulation and many key emergent property statistics (e.g. burst duration and phase gap) are not  
461 calculable or easily approximated with differentiable functions. In such cases, SNPE, which does  
462 not require differentiability of the emergent property has proved to be a powerful approach [35].

463 **Acknowledgements:**

464 This work was funded by NSF Graduate Research Fellowship, DGE-1644869, McKnight Endow-  
465 ment Fund, NIH NINDS 5R01NS100066, Simons Foundation 542963, NSF NeuroNex Award, DBI-  
466 1707398, The Gatsby Charitable Foundation, Simons Collaboration on the Global Brain Postdoc-

467 toral Fellowship, Chinese Postdoctoral Science Foundation, and International Exchange Program  
468 Fellowship. Helpful conversations were had with Francesca Mastrogiuseppe, Srdjan Ostojic, James  
469 Fitzgerald, Stephen Baccus, Dhruva Raman, Liam Paninski, and Larry Abbott.

470 **Data availability statement:**

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

473 **Code availability statement:**

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

475 **References**

- 476 [1] Nancy Kopell and G Bard Ermentrout. Coupled oscillators and the design of central pattern  
477 generators. *Mathematical biosciences*, 90(1-2):87–109, 1988.
- 478 [2] Eve Marder. From biophysics to models of network function. *Annual review of neuroscience*,  
479 21(1):25–45, 1998.
- 480 [3] Larry F Abbott. Theoretical neuroscience rising. *Neuron*, 60(3):489–495, 2008.
- 481 [4] Xiao-Jing Wang. Neurophysiological and computational principles of cortical rhythms in cog-  
482 nition. *Physiological reviews*, 90(3):1195–1268, 2010.
- 483 [5] Timothy O’Leary, Alexander C Sutton, and Eve Marder. Computational models in the age of  
484 large datasets. *Current opinion in neurobiology*, 32:87–94, 2015.
- 485 [6] Ryan N Gutenkunst, Joshua J Waterfall, Fergal P Casey, Kevin S Brown, Christopher R  
486 Myers, and James P Sethna. Universally sloppy parameter sensitivities in systems biology  
487 models. *PLoS Comput Biol*, 3(10):e189, 2007.
- 488 [7] Kamil Erguler and Michael PH Stumpf. Practical limits for reverse engineering of dynamical  
489 systems: a statistical analysis of sensitivity and parameter inferability in systems biology  
490 models. *Molecular BioSystems*, 7(5):1593–1602, 2011.
- 491 [8] Brian K Mannakee, Aaron P Ragsdale, Mark K Transtrum, and Ryan N Gutenkunst. Sloppi-  
492 ness and the geometry of parameter space. In *Uncertainty in Biology*, pages 271–299. Springer,  
493 2016.

- 494 [9] John J Hopfield. Neural networks and physical systems with emergent collective computational  
495 abilities. *Proceedings of the national academy of sciences*, 79(8):2554–2558, 1982.
- 496 [10] Haim Sompolinsky, Andrea Crisanti, and Hans-Jurgen Sommers. Chaos in random neural  
497 networks. *Physical review letters*, 61(3):259, 1988.
- 498 [11] Andrey V Olypher and Ronald L Calabrese. Using constraints on neuronal activity to reveal  
499 compensatory changes in neuronal parameters. *Journal of Neurophysiology*, 98(6):3749–3758,  
500 2007.
- 501 [12] Misha V Tsodyks, William E Skaggs, Terrence J Sejnowski, and Bruce L McNaughton. Para-  
502 doxical effects of external modulation of inhibitory interneurons. *Journal of neuroscience*,  
503 17(11):4382–4388, 1997.
- 504 [13] Kong-Fatt Wong and Xiao-Jing Wang. A recurrent network mechanism of time integration in  
505 perceptual decisions. *Journal of Neuroscience*, 26(4):1314–1328, 2006.
- 506 [14] WR Foster, LH Ungar, and JS Schwaber. Significance of conductances in hodgkin-huxley  
507 models. *Journal of neurophysiology*, 70(6):2502–2518, 1993.
- 508 [15] Astrid A Prinz, Dirk Bucher, and Eve Marder. Similar network activity from disparate circuit  
509 parameters. *Nature neuroscience*, 7(12):1345–1352, 2004.
- 510 [16] Pablo Achard and Erik De Schutter. Complex parameter landscape for a complex neuron  
511 model. *PLoS computational biology*, 2(7):e94, 2006.
- 512 [17] Dmitry Fisher, Itsaso Olasagasti, David W Tank, Emre RF Aksay, and Mark S Goldman.  
513 A modeling framework for deriving the structural and functional architecture of a short-term  
514 memory microcircuit. *Neuron*, 79(5):987–1000, 2013.
- 515 [18] Timothy O’Leary, Alex H Williams, Alessio Franci, and Eve Marder. Cell types, network  
516 homeostasis, and pathological compensation from a biologically plausible ion channel expres-  
517 sion model. *Neuron*, 82(4):809–821, 2014.
- 518 [19] Leandro M Alonso and Eve Marder. Visualization of currents in neural models with similar  
519 behavior and different conductance densities. *Elife*, 8:e42722, 2019.
- 520 [20] Liam Paninski and John P Cunningham. Neural data science: accelerating the experiment-  
521 analysis-theory cycle in large-scale neuroscience. *Current opinion in neurobiology*, 50:232–241,  
522 2018.

- 523 [21] Cristopher M Niell and Michael P Stryker. Modulation of visual responses by behavioral state  
524 in mouse visual cortex. *Neuron*, 65(4):472–479, 2010.
- 525 [22] Aman B Saleem, Asli Ayaz, Kathryn J Jeffery, Kenneth D Harris, and Matteo Carandini.  
526 Integration of visual motion and locomotion in mouse visual cortex. *Nature neuroscience*,  
527 16(12):1864–1869, 2013.
- 528 [23] Simon Musall, Matthew T Kaufman, Ashley L Juavinett, Steven Gluf, and Anne K Church-  
529 land. Single-trial neural dynamics are dominated by richly varied movements. *Nature neuro-  
530 science*, 22(10):1677–1686, 2019.
- 531 [24] Mark A Beaumont, Wenyang Zhang, and David J Balding. Approximate bayesian computation  
532 in population genetics. *Genetics*, 162(4):2025–2035, 2002.
- 533 [25] Paul Marjoram, John Molitor, Vincent Plagnol, and Simon Tavaré. Markov chain monte carlo  
534 without likelihoods. *Proceedings of the National Academy of Sciences*, 100(26):15324–15328,  
535 2003.
- 536 [26] Scott A Sisson, Yanan Fan, and Mark M Tanaka. Sequential monte carlo without likelihoods.  
537 *Proceedings of the National Academy of Sciences*, 104(6):1760–1765, 2007.
- 538 [27] Andreas Raue, Clemens Kreutz, Thomas Maiwald, Julie Bachmann, Marcel Schilling, Ursula  
539 Klingmüller, and Jens Timmer. Structural and practical identifiability analysis of partially  
540 observed dynamical models by exploiting the profile likelihood. *Bioinformatics*, 25(15):1923–  
541 1929, 2009.
- 542 [28] Johan Karlsson, Milena Anguelova, and Mats Jirstrand. An efficient method for structural  
543 identifiability analysis of large dynamic systems. *IFAC Proceedings Volumes*, 45(16):941–946,  
544 2012.
- 545 [29] Keegan E Hines, Thomas R Middendorf, and Richard W Aldrich. Determination of parameter  
546 identifiability in nonlinear biophysical models: A bayesian approach. *Journal of General  
547 Physiology*, 143(3):401–416, 2014.
- 548 [30] Dhruva V Raman, James Anderson, and Antonis Papachristodoulou. Delineating parameter  
549 unidentifiabilities in complex models. *Physical Review E*, 95(3):032314, 2017.
- 550 [31] Gamaleldin F Elsayed and John P Cunningham. Structure in neural population recordings:  
551 an expected byproduct of simpler phenomena? *Nature neuroscience*, 20(9):1310, 2017.

- 552 [32] Cristina Savin and Gašper Tkačik. Maximum entropy models as a tool for building precise  
553 neural controls. *Current opinion in neurobiology*, 46:120–126, 2017.
- 554 [33] Wiktor Mlynarski, Michal Hledík, Thomas R Sokolowski, and Gašper Tkačik. Statistical  
555 analysis and optimality of neural systems. *bioRxiv*, page 848374, 2020.
- 556 [34] Dustin Tran, Rajesh Ranganath, and David Blei. Hierarchical implicit models and likelihood-  
557 free variational inference. In *Advances in Neural Information Processing Systems*, pages 5523–  
558 5533, 2017.
- 559 [35] Pedro J Gonçalves, Jan-Matthis Lueckmann, Michael Deistler, Marcel Nonnenmacher, Kaan  
560 Öcal, Giacomo Bassetto, Chaitanya Chintaluri, William F Podlaski, Sara A Haddad, Tim P  
561 Vogels, et al. Training deep neural density estimators to identify mechanistic models of neural  
562 dynamics. *bioRxiv*, page 838383, 2019.
- 563 [36] Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows.  
564 *International Conference on Machine Learning*, 2015.
- 565 [37] George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji  
566 Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *arXiv preprint*  
567 *arXiv:1912.02762*, 2019.
- 568 [38] Gabriel Loaiza-Ganem, Yuanjun Gao, and John P Cunningham. Maximum entropy flow  
569 networks. *International Conference on Learning Representations*, 2017.
- 570 [39] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp.  
571 *Proceedings of the 5th International Conference on Learning Representations*, 2017.
- 572 [40] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolu-  
573 tions. In *Advances in neural information processing systems*, pages 10215–10224, 2018.
- 574 [41] Gabrielle J Gutierrez, Timothy O’Leary, and Eve Marder. Multiple mechanisms switch an  
575 electrically coupled, synaptically inhibited neuron between competing rhythmic oscillators.  
576 *Neuron*, 77(5):845–858, 2013.
- 577 [42] Mark S Goldman, Jorge Golowasch, Eve Marder, and LF Abbott. Global structure, robustness,  
578 and modulation of neuronal models. *Journal of Neuroscience*, 21(14):5229–5238, 2001.
- 579 [43] Brendan K Murphy and Kenneth D Miller. Balanced amplification: a new mechanism of  
580 selective amplification of neural activity patterns. *Neuron*, 61(4):635–648, 2009.

- 581 [44] Guillaume Hennequin, Tim P Vogels, and Wulfram Gerstner. Optimal control of transient dy-  
582 namics in balanced networks supports generation of complex movements. *Neuron*, 82(6):1394–  
583 1406, 2014.
- 584 [45] Giulio Bondanelli, Thomas Deneux, Brice Bathellier, and Srdjan Ostojic. Population coding  
585 and network dynamics during off responses in auditory cortex. *BioRxiv*, page 810655, 2019.
- 586 [46] Ashok Litwin-Kumar, Robert Rosenbaum, and Brent Doiron. Inhibitory stabilization and vi-  
587 sual coding in cortical circuits with multiple interneuron subtypes. *Journal of neurophysiology*,  
588 115(3):1399–1409, 2016.
- 589 [47] Agostina Palmigiano, Francesco Fumarola, Daniel P Mossing, Nataliya Kraynyukova, Hillel  
590 Adesnik, and Kenneth Miller. Structure and variability of optogenetic responses identify the  
591 operating regime of cortex. *bioRxiv*, 2020.
- 592 [48] Chunyu A Duan, Marino Pagan, Alex T Piet, Charles D Kopec, Athena Akrami, Alexander J  
593 Riordan, Jeffrey C Erlich, and Carlos D Brody. Collicular circuits for flexible sensorimotor  
594 routing. *bioRxiv*, page 245613, 2018.
- 595 [49] Eve Marder and Vatsala Thirumalai. Cellular, synaptic and network effects of neuromodula-  
596 tion. *Neural Networks*, 15(4-6):479–493, 2002.
- 597 [50] Mark S Goldman. Memory without feedback in a neural network. *Neuron*, 61(4):621–634,  
598 2009.
- 599 [51] Giulio Bondanelli and Srdjan Ostojic. Coding with transient trajectories in recurrent neural  
600 networks. *PLoS computational biology*, 16(2):e1007655, 2020.
- 601 [52] David Sussillo. Neural circuits as computational dynamical systems. *Current opinion in*  
602 *neurobiology*, 25:156–163, 2014.
- 603 [53] Omri Barak. Recurrent neural networks as versatile tools of neuroscience research. *Current*  
604 *opinion in neurobiology*, 46:1–6, 2017.
- 605 [54] Abigail A Russo, Sean R Bittner, Sean M Perkins, Jeffrey S Seely, Brian M London, Antonio H  
606 Lara, Andrew Miri, Najja J Marshall, Adam Kohn, Thomas M Jessell, et al. Motor cortex  
607 embeds muscle-like commands in an untangled population response. *Neuron*, 97(4):953–966,  
608 2018.

- 609 [55] Scott A Sisson, Yanan Fan, and Mark Beaumont. *Handbook of approximate Bayesian computation*. CRC Press, 2018.
- 610
- 611 [56] Kyle Cranmer, Johann Brehmer, and Gilles Louppe. The frontier of simulation-based inference.
- 612 *Proceedings of the National Academy of Sciences*, 2020.
- 613 [57] Eve Marder and Allen I Selverston. *Dynamic biological networks: the stomatogastric nervous system*. MIT press, 1992.
- 614
- 615 [58] Hirofumi Ozeki, Ian M Finn, Evan S Schaffer, Kenneth D Miller, and David Ferster. Inhibitory stabilization of the cortical network underlies visual surround suppression. *Neuron*, 62(4):578–592, 2009.
- 616
- 617
- 618 [59] Daniel B Rubin, Stephen D Van Hooser, and Kenneth D Miller. The stabilized supralinear network: a unifying circuit motif underlying multi-input integration in sensory cortex. *Neuron*, 85(2):402–417, 2015.
- 619
- 620
- 621 [60] Guillaume Hennequin, Yashar Ahmadian, Daniel B Rubin, Máté Lengyel, and Kenneth D Miller. The dynamical regime of sensory cortex: stable dynamics around a single stimulus-tuned attractor account for patterns of noise variability. *Neuron*, 98(4):846–860, 2018.
- 622
- 623
- 624 [61] Mark M. Churchland, Byron M. Yu, John P. Cunningham, Leo P. Sugrue, Marlene R. Cohen, Greg S. Corrado, William T. Newsome, Andrew M. Clark, Paymon Hosseini, Benjamin B. Scott, David C. Bradley, Matthew A. Smith, Adam Kohn, J. Anthony Movshon, Katherine M. Armstrong, Tirin Moore, Steve W. Chang, Lawrence H. Snyder, Stephen G. Lisberger, Nicholas J. Priebe, Ian M. Finn, David Ferster, Stephen I. Ryu, Gopal Santhanam, Maneesh Sahani, and Krishna V. Shenoy. Stimulus onset quenches neural variability: a widespread cortical phenomenon. *Nat. Neurosci.*, 13(3):369–378, 2010.
- 625
- 626
- 627
- 628
- 629
- 630
- 631 [62] João D Semedo, Amin Zandvakili, Christian K Machens, M Yu Byron, and Adam Kohn. Cortical areas interact through a communication subspace. *Neuron*, 102(1):249–259, 2019.
- 632
- 633 [63] Henry Markram, Maria Toledo-Rodriguez, Yun Wang, Anirudh Gupta, Gilad Silberberg, and Caizhi Wu. Interneurons of the neocortical inhibitory system. *Nature reviews neuroscience*, 5(10):793, 2004.
- 634
- 635
- 636 [64] Bernardo Rudy, Gordon Fishell, SooHyun Lee, and Jens Hjerling-Leffler. Three groups of interneurons account for nearly 100% of neocortical gabaergic neurons. *Developmental neurobiology*, 71(1):45–61, 2011.
- 637
- 638

- 639 [65] Robin Tremblay, Soohyun Lee, and Bernardo Rudy. GABAergic Interneurons in the Neocortex:  
640 From Cellular Properties to Circuits. *Neuron*, 91(2):260–292, 2016.
- 641 [66] Carsten K Pfeffer, Mingshan Xue, Miao He, Z Josh Huang, and Massimo Scanziani. Inhi-  
642 bition of inhibition in visual cortex: the logic of connections between molecularly distinct  
643 interneurons. *Nature Neuroscience*, 16(8):1068, 2013.
- 644 [67] Daniel J Felleman and David C Van Essen. Distributed hierarchical processing in the primate  
645 cerebral cortex. *Cerebral cortex (New York, NY: 1991)*, 1(1):1–47, 1991.
- 646 [68] C Gardiner. Stochastic methods: A Handbook for the Natural and Social Sciences, 2009.
- 647 [69] Chunyu A Duan, Jeffrey C Erlich, and Carlos D Brody. Requirement of prefrontal and midbrain  
648 regions for rapid executive control of behavior in the rat. *Neuron*, 86(6):1491–1503, 2015.
- 649 [70] Mark Girolami and Ben Calderhead. Riemann manifold langevin and hamiltonian monte  
650 carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*,  
651 73(2):123–214, 2011.
- 652 [71] Lawrence Saul and Michael Jordan. A mean field learning algorithm for unsupervised neural  
653 networks. In *Learning in graphical models*, pages 541–554. Springer, 1998.
- 654 [72] Nicholas Metropolis, Arianna W Rosenbluth, Marshall N Rosenbluth, Augusta H Teller, and  
655 Edward Teller. Equation of state calculations by fast computing machines. *The journal of  
656 chemical physics*, 21(6):1087–1092, 1953.
- 657 [73] W Keith Hastings. Monte carlo sampling methods using markov chains and their applications.  
658 1970.
- 659 [74] Ben Calderhead and Mark Girolami. Statistical analysis of nonlinear dynamical systems using  
660 differential geometric sampling methods. *Interface focus*, 1(6):821–835, 2011.
- 661 [75] Andrew Golightly and Darren J Wilkinson. Bayesian parameter inference for stochastic bio-  
662 chemical network models using particle markov chain monte carlo. *Interface focus*, 1(6):807–  
663 820, 2011.
- 664 [76] Oksana A Chkrebtii, David A Campbell, Ben Calderhead, Mark A Girolami, et al. Bayesian  
665 solution uncertainty quantification for differential equations. *Bayesian Analysis*, 11(4):1239–  
666 1267, 2016.

- 667 [77] Juliane Liepe, Paul Kirk, Sarah Filippi, Tina Toni, Chris P Barnes, and Michael PH Stumpf.  
668 A framework for parameter estimation and model selection from experimental data in systems  
669 biology using approximate bayesian computation. *Nature protocols*, 9(2):439–456, 2014.
- 670 [78] Sean R Bittner, Agostina Palmigiano, Kenneth D Miller, and John P Cunningham. Degener-  
671 ate solution networks for theoretical neuroscience. *Computational and Systems Neuroscience*  
672 *Meeting (COSYNE), Lisbon, Portugal*, 2019.
- 673 [79] Sean R Bittner, Alex T Piet, Chunyu A Duan, Agostina Palmigiano, Kenneth D Miller,  
674 Carlos D Brody, and John P Cunningham. Examining models in theoretical neuroscience with  
675 degenerate solution networks. *Bernstein Conference 2019, Berlin, Germany*, 2019.
- 676 [80] Marcel Nonnenmacher, Pedro J Goncalves, Giacomo Bassetto, Jan-Matthis Lueckmann, and  
677 Jakob H Macke. Robust statistical inference for simulation-based models in neuroscience. In  
678 *Bernstein Conference 2018, Berlin, Germany*, 2018.
- 679 [81] Deistler Michael, , Pedro J Goncalves, Kaan Oecal, and Jakob H Macke. Statistical inference for  
680 analyzing sloppiness in neuroscience models. In *Bernstein Conference 2019, Berlin, Germany*,  
681 2019.
- 682 [82] Jan-Matthis Lueckmann, Pedro J Goncalves, Giacomo Bassetto, Kaan Öcal, Marcel Nonnen-  
683 macher, and Jakob H Macke. Flexible statistical inference for mechanistic models of neural  
684 dynamics. In *Advances in Neural Information Processing Systems*, pages 1289–1299, 2017.
- 685 [83] George Papamakarios, David Sterratt, and Iain Murray. Sequential neural likelihood: Fast  
686 likelihood-free inference with autoregressive flows. In *The 22nd International Conference on*  
687 *Artificial Intelligence and Statistics*, pages 837–848. PMLR, 2019.
- 688 [84] Joeri Hermans, Volodimir Begy, and Gilles Louppe. Likelihood-free mcmc with amortized  
689 approximate ratio estimators. In *International Conference on Machine Learning*, pages 4239–  
690 4248. PMLR, 2020.
- 691 [85] Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and  
692 variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.
- 693 [86] Sean R Bittner and John P Cunningham. Approximating exponential family models (not  
694 single distributions) with a two-network architecture. *arXiv preprint arXiv:1903.07515*, 2019.

- 695 [87] Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary  
696 differential equations. In *Advances in neural information processing systems*, pages 6571–6583,  
697 2018.
- 698 [88] Xuechen Li, Ting-Kam Leonard Wong, Ricky TQ Chen, and David Duvenaud. Scalable  
699 gradients for stochastic differential equations. *arXiv preprint arXiv:2001.01328*, 2020.
- 700 [89] Maria Pia Saccomani, Stefania Audoly, and Leontina D’Angiò. Parameter identifiability of  
701 nonlinear systems: the role of initial conditions. *Automatica*, 39(4):619–632, 2003.
- 702 [90] Stefan Hengl, Clemens Kreutz, Jens Timmer, and Thomas Maiwald. Data-based identifiability  
703 analysis of non-linear dynamical models. *Bioinformatics*, 23(19):2612–2618, 2007.
- 704 [91] George Papamakarios, Theo Pavlakou, and Iain Murray. Masked autoregressive flow for density  
705 estimation. In *Advances in Neural Information Processing Systems*, pages 2338–2347, 2017.
- 706 [92] Durk P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling.  
707 Improved variational inference with inverse autoregressive flow. *Advances in neural information  
708 processing systems*, 29:4743–4751, 2016.
- 709 [93] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *International  
710 Conference on Learning Representations*, 2015.
- 711 [94] Emmanuel Klinger, Dennis Rickert, and Jan Hasenauer. pyabc: distributed, likelihood-free  
712 inference. *Bioinformatics*, 34(20):3591–3593, 2018.
- 713 [95] David S Greenberg, Marcel Nonnenmacher, and Jakob H Macke. Automatic posterior trans-  
714 formation for likelihood-free inference. *International Conference on Machine Learning*, 2019.
- 715 [96] Daniel P Mossing, Julia Veit, Agostina Palmigiano, Kenneth D. Miller, and Hillel Adesnik.  
716 Antagonistic inhibitory subnetworks control cooperation and competition across cortical space.  
717 *bioRxiv*, 2021.

718 **5 Methods**

719 **5.1 Emergent property inference (EPI)**

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

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

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

$$\mathcal{X} : \mathbb{E}_{\mathbf{z}, \mathbf{x}} [f(\mathbf{x}; \mathbf{z})] = \boldsymbol{\mu}, \text{Var}_{\mathbf{z}, \mathbf{x}} [f(\mathbf{x}; \mathbf{z})] = \boldsymbol{\sigma}^2. \quad (10)$$

741 Precisely, the emergent property statistics  $f(\mathbf{x}; \mathbf{z})$  must have means  $\boldsymbol{\mu}$  and variances  $\boldsymbol{\sigma}^2$  over the EPI  
742 distribution of parameters and the data produced by those parameters. Technically, an emergent  
743 property may be a combination of first-, second-, or higher-order moments, but this study focuses  
744 on the case written in Equation 10.

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

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

$$\mathbf{z} = g_{\theta}(\mathbf{z}_0) = g_l(\dots g_1(\mathbf{z}_0)) \sim q_{\theta}(\mathbf{z}). \quad (11)$$

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

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

$$q_{\theta}(\mathbf{z} | \mathcal{X}) = q_{\theta}^*(\mathbf{z}) = \underset{q_{\theta} \in \mathcal{Q}}{\operatorname{argmax}} H(q_{\theta}(\mathbf{z})) \quad (12)$$

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

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

760 In the remainder of Section 5.1, we will explain the finer details and motivation of the EPI method.  
 761 First, we explain related approaches and what EPI introduces to this domain (Section 5.1.1). Sec-  
 762 ond, we describe the special class of deep probability distributions used in EPI called normalizing  
 763 flows (Section 5.1.2). Next, we explain the constrained optimization technique used to solve Equa-  
 764 tion 12 (Section 5.1.4). Then, we demonstrate the details of this optimization in a toy example  
 765 (Section 5.1.5). Finally, we establish the known relationship between maximum entropy distri-  
 766 butions and exponential families (Section 5.1.3), which is used to explain how EPI is a form of  
 767 variational inference (Section 5.1.6).

### 768 5.1.1 Related approaches

769 When bayesian inference problems lack conjugacy, scientists use approximate inference methods like  
 770 variational inference (VI) [71] and Markov chain Monte Carlo (MCMC) [72, 73]. After optimization,  
 771 variational methods return a parameterized posterior distribution, which we can analyze. Also, the  
 772 variational approximating distribution class is often chosen such that it permits fast sampling. In

773 contrast MCMC methods only produce samples from the approximated posterior distribution. No  
774 parameterized distribution is estimated, and additional samples are always generated with the same  
775 sampling complexity. Inference in models defined by systems of differential has been demonstrated  
776 with MCMC [70], although this approach requires tractable likelihoods. Advancements have intro-  
777 duced sampling [74], likelihood approximation [75], and uncertainty quantification techniques [76]  
778 to make MCMC approaches more efficient and expand the class of applicable models.

779 Simulation-based inference [56] is model parameter inference in the absence of a tractable likeli-  
780 hood function. The most prevalent approach to simulation-based inference is approximate Bayesian  
781 computation (ABC) [24], in which satisfactory parameter samples are kept from random prior sam-  
782 pling according to a rejection heuristic. The obtained set of parameters do not have a probabilities,  
783 and further insight about the model must be gained from examination of the parameter set and  
784 their generated activity. Methodological advances to ABC methods have come through the use of  
785 Markov chain Monte Carlo (MCMC-ABC) [25] and sequential Monte Carlo (SMC-ABC) [26] sam-  
786 pling techniques. SMC-ABC is considered state-of-the-art ABC, yet this approach still struggles  
787 to scale in dimensionality [55] (cf. Fig. 2). Still, this method has enjoyed much success in systems  
788 biology [77]. Furthermore, once a parameter set has been obtained by SMC-ABC from a finite set  
789 of particles, the SMC-ABC algorithm must be run again from scratch with a new population of  
790 initialized particles to obtain additional samples.

791 For scientific model analysis, we seek a parameter distribution represented by an approximating  
792 distribution as in variational inference [71]: a variational approximation that once optimized yields  
793 fast analytic calculations and samples. For the reasons described above, ABC and MCMC tech-  
794 niques are unattractive, since they only produce a set of parameter samples lacking probabilities  
795 and have unchanging sampling rate. EPI infers parameters in circuit models using the MEFN [38]  
796 algorithm with a deep variational approximation. The deep neural network of EPI (Fig. 1E) de-  
797 fines the parametric form (with variational parameters  $\theta$ ) of the deep variational approximation of  
798 circuit parameters  $\mathbf{z}$ .

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

804 We note that, during our preparation and early presentation of this work [78, 79], another work

805 has arisen with broadly similar goals: bringing statistical inference to mechanistic models of neural  
806 circuits ([35, 80, 81]). We are encouraged by this general problem being recognized by others in the  
807 community, and we emphasize that these works offer complementary neuroscientific contributions  
808 (different theoretical models of focus) and use different technical methodologies (ours is built on  
809 our prior work [38], theirs similarly [82]).

810 The method EPI differs from SNPE in some key ways. SNPE belongs to a “sequential” class  
811 of recently developed simulation-based inference methods in which two neural networks are used  
812 for posterior inference. This first neural network is a deep probability distribution (normalizing  
813 flow) used to estimate the posterior  $p(\mathbf{z} | \mathbf{x})$  (SNPE) or the likelihood  $p(\mathbf{x} | \mathbf{z})$  (sequential neural  
814 likelihood (SNL [83])). A recent advance uses an unconstrained neural network to estimate the  
815 likelihood ratio (sequential neural ratio estimation (SNRE [84])). In SNL and SNRE, MCMC  
816 sampling techniques are used to obtain samples from the approximated posterior. This contrasts  
817 with EPI and SNPE, which use deep probability distributions to model parameters, which facilitates  
818 immediate measurements of sample probability, gradient, or Hessian for system analysis. The  
819 second neural network in this sequential class of methods is the amortizer. This unconstrained  
820 deep network maps data  $\mathbf{x}$  (or statistics  $f(\mathbf{x}; \mathbf{z})$ ) or model parameters  $\mathbf{z}$  to the weights and biases of  
821 the first neural network. These methods are optimized on a conditional density (or ratio) estimation  
822 objective. The data used to optimize this objective are generated via an adaptive procedure, in  
823 which training data pairs  $(\mathbf{x}_i, \mathbf{z}_i)$  become sequentially closer to the true data and posterior.

824 The approximating fidelity of the deep probability distribution in sequential approaches is opti-  
825 mized to generalize across the training distribution of the conditioning variable. This generalization  
826 property of the sequential methods can reduce the accuracy at the singular posterior of interest.  
827 Whereas in EPI, the entire expressivity of the deep probability distribution is dedicated to learning  
828 a single distribution as well as possible. Amortization is not possible in EPI, since EPI learns  
829 an exponential family distribution parameterized by its mean (see Section 5.1.3). Since EPI dis-  
830 tributions are defined by the mean  $\mu$  of their statistics, there is the well-known inverse mapping  
831 problem of exponential families [85] that prohibits an amortization based approach. However, we  
832 have shown that the same two-network architecture of the sequential simulation-based inference  
833 methods can be used for amortized inference in intractable exponential family posteriors using their  
834 natural parameterization [86].

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

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

845 Structural identifiability analysis involves the measurement of sensitivity and unidentifiabilities in  
 846 scientific models. Around a single parameter choice, one can measure the Jacobian. One approach  
 847 for this calculation that scales well is EAR [28]. A popular efficient approach for systems of ODEs  
 848 has been neural ODE adjoint [87] and its stochastic adaptation [88]. Casting identifiability as a  
 849 statistical estimation problem, the profile likelihood works via iterated optimization while holding  
 850 parameters fixed [27]. An exciting recent method is capable of recovering the functional form of such  
 851 unidentifiabilities away from a point by following degenerate dimensions of the fisher information  
 852 matrix [30]. Global structural non-identifiabilities can be found for models with polynomial or  
 853 rational dynamics equations using DAISY [89], or through mean optimal transformations [90].  
 854 With EPI, we have all the benefits given by a statistical inference method plus the ability to query  
 855 the first- or second-order gradient of the probability of the inferred distribution at any chosen  
 856 parameter value. The second-order gradient of the log probability (the Hessian), which is directly  
 857 afforded by EPI distributions, produces quantified information about parametric sensitivity of the  
 858 emergent property in parameter space (see Section 3.2).

### 859 **5.1.2 Deep probability distributions and normalizing flows**

860 Deep probability distributions are comprised of multiple layers of fully connected neural networks  
 861 (Equation 11). When each neural network layer is restricted to be a bijective function, the sample  
 862 density can be calculated using the change of variables formula at each layer of the network. For  
 863  $\mathbf{z}_i = g_i(\mathbf{z}_{i-1})$ ,

$$p(\mathbf{z}_i) = p(g_i^{-1}(\mathbf{z}_i)) \left| \det \frac{\partial g_i^{-1}(\mathbf{z}_i)}{\partial \mathbf{z}_i} \right| = p(\mathbf{z}_{i-1}) \left| \det \frac{\partial g_i(\mathbf{z}_{i-1})}{\partial \mathbf{z}_{i-1}} \right|^{-1}. \quad (13)$$

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

866 Jacobian computations, which confer a fast calculation of the sample log probability. Fast log  
867 probability calculation confers efficient optimization of the maximum entropy objective (see Section  
868 5.1.4).

869 We use the Real NVP [39] normalizing flow class, because its coupling architecture confers both  
870 fast sampling (forward) and fast log probability evaluation (backward). Fast probability evaluation  
871 facilitates fast gradient and Hessian evaluation of log probability throughout parameter space.  
872 Glow permutations were used in between coupling stages [40]. This is in contrast to autoregressive  
873 architectures [91, 92], in which only one of the forward or backward passes can be efficient. In this  
874 work, normalizing flows are used as flexible parameter distribution approximations  $q_{\theta}(\mathbf{z})$  having  
875 weights and biases  $\theta$ . We specify the architecture used in each application by the number of Real-  
876 NVP affine coupling stages, and the number of neural network layers and units per layer of the  
877 conditioning functions.

878 When calculating Hessians of log probabilities in deep probability distributions, it is important to  
879 consider the normalizing flow architecture. With autoregressive architectures [91, 92], fast sam-  
880 pling and fast log probability evaluations are mutually exclusive. That makes these architectures  
881 undesirable for EPI, where efficient sampling is important for optimization, and log probability  
882 evaluation speed predicates the efficiency of gradient and Hessian calculations. With Real NVP  
883 coupling architectures, we get both fast sampling and fast Hessians making both optimization and  
884 scientific analysis efficient.

### 885 5.1.3 Maximum entropy distributions and exponential families

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

$$p^*(\mathbf{z}) = \underset{p \in \mathcal{P}}{\operatorname{argmax}} H(p(\mathbf{z})) \quad (14)$$

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

888 will have probability density in the exponential family:

$$p^*(\mathbf{z}) \propto \exp(\boldsymbol{\eta}^\top T(\mathbf{z})). \quad (15)$$

889 The mappings between the mean parameterization  $\boldsymbol{\mu}_{\text{opt}}$  and the natural parameterization  $\boldsymbol{\eta}$  are  
890 formally hard to identify except in special cases [85].

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

$$\text{Var}_{\mathbf{z}, \mathbf{x}} [f(\mathbf{x}; \mathbf{z})] = \mathbb{E}_{\mathbf{z}, \mathbf{x}} [(f(\mathbf{x}; \mathbf{z}) - \boldsymbol{\mu})^2] \quad (16)$$

893 As a general maximum entropy distribution (Equation 14), the sufficient statistics vector contains  
 894 both first and second order moments of  $f(\mathbf{x}; \mathbf{z})$

$$T(\mathbf{x}; \mathbf{z}) = \begin{bmatrix} f(\mathbf{x}; \mathbf{z}) \\ (f(\mathbf{x}; \mathbf{z}) - \boldsymbol{\mu})^2 \end{bmatrix}, \quad (17)$$

895 which are constrained to the chosen means and variances

$$\boldsymbol{\mu}_{\text{opt}} = \begin{bmatrix} \boldsymbol{\mu} \\ \sigma^2 \end{bmatrix}. \quad (18)$$

#### 896 5.1.4 Augmented Lagrangian optimization

897 To optimize  $q_{\boldsymbol{\theta}}(\mathbf{z})$  in Equation 12, the constrained maximum entropy optimization is executed using  
 898 the augmented Lagrangian method. The following objective is minimized:

$$L(\boldsymbol{\theta}; \boldsymbol{\eta}_{\text{opt}}, c) = -H(q_{\boldsymbol{\theta}}) + \boldsymbol{\eta}_{\text{opt}}^\top R(\boldsymbol{\theta}) + \frac{c}{2} \|R(\boldsymbol{\theta})\|^2 \quad (19)$$

899 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  
 900 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  
 901 sufficient statistics  $T(\mathbf{x}; \mathbf{z})$  and mean parameter  $\boldsymbol{\mu}_{\text{opt}}$  are determined by the means  $\boldsymbol{\mu}$  and variances  
 902  $\sigma^2$  of emergent property statistics  $f(\mathbf{x}; \mathbf{z})$  defined in Equation 12 (see Section 5.1.6). Specifically,  
 903  $T(\mathbf{x}; \mathbf{z})$  is a concatenation of the first and second moments,  $\boldsymbol{\mu}_{\text{opt}}$  is a concatenation of  $\boldsymbol{\mu}$  and  $\sigma^2$   
 904 (see section 5.1.3), and the Lagrange multipliers are closely related to the natural parameters  $\boldsymbol{\eta}$  of  
 905 exponential families (see Section 5.1.3). Weights and biases  $\boldsymbol{\theta}$  of the deep probability distribution  
 906 are optimized according to Equation 19 using the Adam optimizer with learning rate  $10^{-3}$  [93].

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

$$H(q_{\boldsymbol{\theta}}(\mathbf{z})) = \int -q_{\boldsymbol{\theta}}(\mathbf{z}) \log(q_{\boldsymbol{\theta}}(\mathbf{z})) d\mathbf{z} = \mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\theta}}} [-\log(q_{\boldsymbol{\theta}}(\mathbf{z}))] = \mathbb{E}_{\mathbf{z}_0 \sim q_0} [-\log(q_{\boldsymbol{\theta}}(g_{\boldsymbol{\theta}}(\mathbf{z}_0)))]. \quad (20)$$

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

$$\nabla_{\boldsymbol{\theta}} H(q_{\boldsymbol{\theta}}(\mathbf{z})) = \mathbb{E}_{\mathbf{z}_0 \sim q_0} [-\nabla_{\boldsymbol{\theta}} \log(q_{\boldsymbol{\theta}}(g_{\boldsymbol{\theta}}(\mathbf{z}_0)))]. \quad (21)$$

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

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

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

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

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

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

### 954 5.1.5 Example: 2D LDS

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

$$957 \quad \tau \frac{d\mathbf{x}}{dt} = A\mathbf{x} \quad (22)$$

957 with

$$958 \quad A = \begin{bmatrix} a_1 & a_2 \\ a_3 & a_4 \end{bmatrix}. \quad (23)$$

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

$$966 \quad \mathbb{E}[T(\mathbf{x})]_{\mathbf{z}, \mathbf{x}} \triangleq \mathbb{E} \begin{bmatrix} \text{real}(\lambda_1)(\mathbf{x}; \mathbf{z}) \\ \text{imag}(\lambda_1)(\mathbf{x}; \mathbf{z}) \\ (\text{real}(\lambda_1)(\mathbf{x}; \mathbf{z}) - 0)^2 \\ (\text{imag}(\lambda_1)(\mathbf{x}; \mathbf{z}) - 2\pi\omega)^2 \end{bmatrix} = \begin{bmatrix} 0.0 \\ 2\pi \\ 0.25^2 \\ (2\pi 0.1)^2 \end{bmatrix} \triangleq \boldsymbol{\mu}_{\text{opt}}. \quad (24)$$

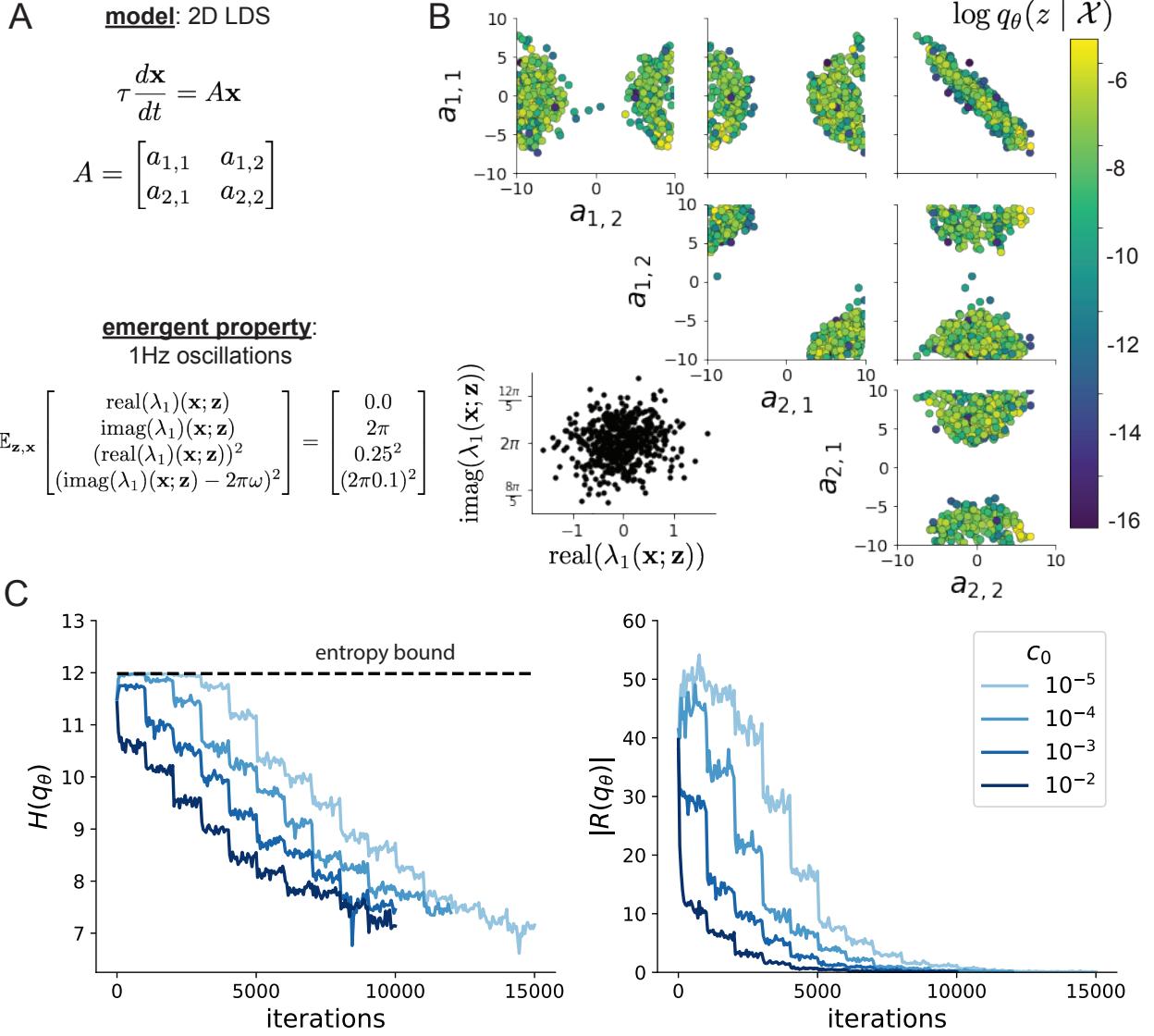


Figure 6: (LDS1): **A.** Two-dimensional linear dynamical system model, where real entries of the dynamics matrix  $A$  are the parameters. **B.** The EPI distribution for a two-dimensional linear dynamical system with  $\tau = 1$  that produces an average of 1Hz oscillations with some small amount of 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.

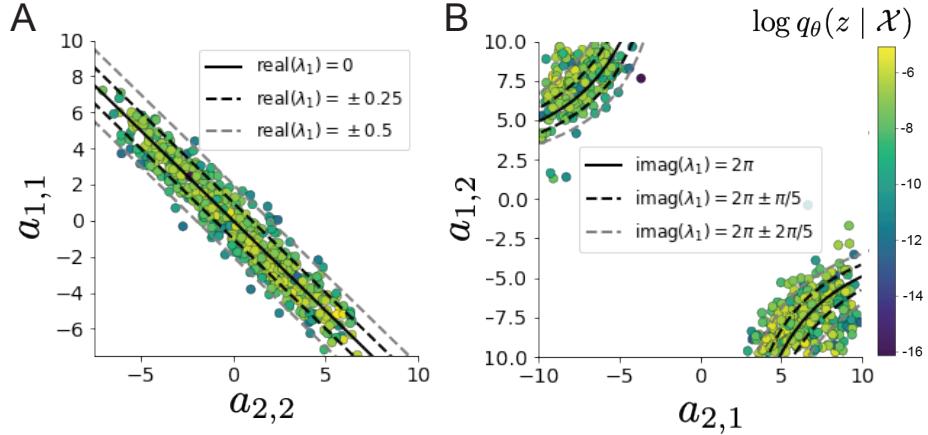


Figure 7: (LDS2): **A.** Probability contours in the  $a_1$ - $a_4$  plane were derived from the relationship to emergent property statistic of growth/decay factor  $\text{real}(\lambda_1)$ . **B.** Probability contours in the  $a_2$ - $a_3$  plane were derived from the emergent property statistic of oscillation frequency  $2\pi\text{imag}(\lambda_1)$ .

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

$$\lambda = \frac{\left(\frac{a_1+a_4}{\tau}\right) \pm \sqrt{\left(\frac{a_1+a_4}{\tau}\right)^2 + 4\left(\frac{a_2a_3-a_1a_4}{\tau}\right)}}{2}. \quad (25)$$

970 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  
 971 does not require simulation, we cannot derive the distribution  $q_\theta^*$  directly. This fact is due to the  
 972 formally hard problem of the backward mapping: finding the natural parameters  $\eta$  from the mean  
 973 parameters  $\mu$  of an exponential family distribution [85]. Instead, we used EPI to approximate this  
 974 distribution (Fig. S1B). We used a real-NVP normalizing flow architecture with four masks, two  
 975 neural network layers of 15 units per mask, with batch normalization momentum 0.99, mapped  
 976 onto a support of  $z_i \in [-10, 10]$ . (see Section 5.1.2).

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

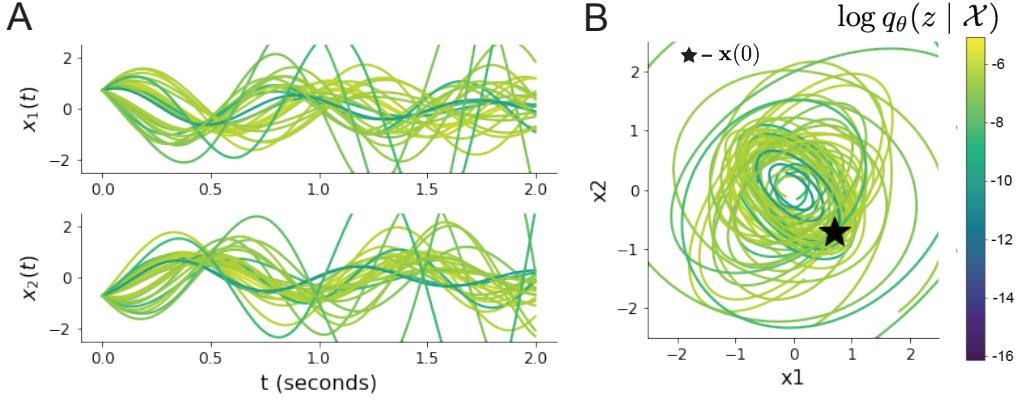


Figure 8: (LDS3): Sampled dynamical systems  $\mathbf{z} \sim q_\theta(\mathbf{z})$  and their simulated activity from  $\mathbf{x}(0) = [\frac{\sqrt{2}}{2}, -\frac{\sqrt{2}}{2}]$  colored by log probability. **A.** Each dimension of the simulated trajectories throughout time. **B.** The simulated trajectories in phase space.

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

$$\text{imag}(\lambda_1) = \begin{cases} \sqrt{\frac{a_1a_4-a_2a_3}{\tau}}, & \text{if } a_1a_4 < a_2a_3 \\ 0 & \text{otherwise} \end{cases}. \quad (26)$$

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

$$\text{imag}(\lambda_1)^2 = a_1a_4 - a_2a_3 \quad (27)$$

989 Since we constrained  $\mathbb{E}_{\mathbf{z} \sim q_\theta} [\text{imag}(\lambda)] = 2\pi$ , we can plot contours of the equation  $\text{imag}(\lambda_1)^2 =$   
 990  $a_1a_4 - a_2a_3 = (2\pi)^2$  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  
 991 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).  
 992 This validates the curved structure of the inferred distribution learned through EPI. We took steps  
 993 in negative standard deviation of  $a_1a_4$  (dotted and gray lines), since there are few positive values  
 994  $a_1a_4$  in the learned distribution. Subtler combinations of model and emergent property will have  
 995 more complexity, further motivating the use of EPI for understanding these systems. As we expect,  
 996 the distribution results in samples of two-dimensional linear systems oscillating near 1Hz (Fig. S3).

997 **5.1.6 EPI as variational inference**

998 In Bayesian inference a prior belief about model parameters  $\mathbf{z}$  is stated in a prior distribution  $p(\mathbf{z})$ ,  
 999 and the statistical model capturing the effect of  $\mathbf{z}$  on observed data points  $\mathbf{x}$  is formalized in the  
 1000 likelihood distribution  $p(\mathbf{x} | \mathbf{z})$ . In Bayesian inference, we obtain a posterior distribution  $p(\mathbf{z} | \mathbf{x})$ ,  
 1001 which captures how the data inform our knowledge of model parameters using Bayes' rule:

$$p(\mathbf{z} | \mathbf{x}) = \frac{p(\mathbf{x} | \mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}. \quad (28)$$

1002 The posterior distribution is analytically available when the prior is conjugate with the likelihood.  
 1003 However, conjugacy is rare in practice, and alternative methods, such as variational inference [71],  
 1004 are utilized.

1005 In variational inference, a posterior approximation  $q_{\theta}^*$  is chosen from within some variational family  
 1006  $\mathcal{Q}$

$$q_{\theta}^*(\mathbf{z}) = \operatorname{argmin}_{q_{\theta} \in \mathcal{Q}} KL(q_{\theta}(\mathbf{z}) || p(\mathbf{z} | \mathbf{x})). \quad (29)$$

1007 The KL divergence can be written in terms of entropy of the variational approximation:

$$KL(q_{\theta}(\mathbf{z}) || p(\mathbf{z} | \mathbf{x})) = \mathbb{E}_{\mathbf{z} \sim q_{\theta}} [\log(q_{\theta}(\mathbf{z}))] - \mathbb{E}_{\mathbf{z} \sim q_{\theta}} [\log(p(\mathbf{z} | \mathbf{x}))] \quad (30)$$

$$= -H(q_{\theta}) - \mathbb{E}_{\mathbf{z} \sim q_{\theta}} [\log(p(\mathbf{x} | \mathbf{z})) + \log(p(\mathbf{z})) - \log(p(\mathbf{x}))] \quad (31)$$

1009 Since the marginal distribution of the data  $p(\mathbf{x})$  (or “evidence”) is independent of  $\theta$ , variational  
 1010 inference is executed by optimizing the remaining expression. This is usually framed as maximizing  
 1011 the evidence lower bound (ELBO)

$$\operatorname{argmin}_{q_{\theta} \in \mathcal{Q}} KL(q_{\theta} || p(\mathbf{z} | \mathbf{x})) = \operatorname{argmax}_{q_{\theta} \in \mathcal{Q}} H(q_{\theta}) + \mathbb{E}_{\mathbf{z} \sim q_{\theta}} [\log(p(\mathbf{x} | \mathbf{z})) + \log(p(\mathbf{z}))]. \quad (32)$$

1012 Now, consider the setting where we have chosen a uniform prior, and stipulate a mean-field gaussian  
 1013 likelihood on a chosen statistic of the data  $f(\mathbf{x}; \mathbf{z})$

$$p(\mathbf{x} | \mathbf{z}) = \mathcal{N}(f(\mathbf{x}; \mathbf{z}) | \boldsymbol{\mu}_f, \Sigma_f), \quad (33)$$

1014 where  $\Sigma_f = \text{diag}(\boldsymbol{\sigma}_f^2)$ . The log likelihood is then proportional to a dot product of the natural  
 1015 parameter of this mean-field gaussian distribution and the first and second moment statistics.

$$\log p(\mathbf{x} | \mathbf{z}) \propto \boldsymbol{\eta}_f^\top T(\mathbf{x}, \mathbf{z}), \quad (34)$$

1016 where

$$\boldsymbol{\eta}_f = \begin{bmatrix} \frac{\boldsymbol{\mu}_f}{\boldsymbol{\sigma}_f^2} \\ \frac{-1}{2\boldsymbol{\sigma}_f^2} \end{bmatrix}, \text{ and} \quad (35)$$

1017

$$T(\mathbf{x}; \mathbf{z}) = \begin{bmatrix} f(\mathbf{x}; \mathbf{z}) \\ (f(\mathbf{x}; \mathbf{z}) - \boldsymbol{\mu}_f)^2 \end{bmatrix}. \quad (36)$$

1018 The variational objective is then

$$\underset{q_{\theta} \in Q}{\operatorname{argmax}} H(q_{\theta}) + \boldsymbol{\eta}_f^{\top} \mathbb{E}_{\mathbf{z} \sim q_{\theta}} [T(\mathbf{x}; \mathbf{z})] \quad (37)$$

1019 Comparing this to the Lagrangian objective (without augmentation) of EPI, we see they are the  
1020 same

$$\begin{aligned} q_{\theta}^*(\mathbf{z}) &= \underset{q_{\theta} \in Q}{\operatorname{argmin}} -H(q_{\theta}) + \boldsymbol{\eta}_{\text{opt}}^{\top} (\mathbb{E}_{\mathbf{z}, \mathbf{x}} [T(\mathbf{x}; \mathbf{z})] - \boldsymbol{\mu}_{\text{opt}}) \\ &= \underset{q_{\theta} \in Q}{\operatorname{argmin}} -H(q_{\theta}) + \boldsymbol{\eta}_{\text{opt}}^{\top} \mathbb{E}_{\mathbf{z}, \mathbf{x}} [T(\mathbf{x}; \mathbf{z})]. \end{aligned} \quad (38)$$

1021 where  $T(\mathbf{x}; \mathbf{z})$  consists of the first and second moments of the emergent property statistic  $f(\mathbf{x}; \mathbf{z})$   
1022 (Equation 17). Thus, EPI is implicitly executing variational inference with a uniform prior and a  
1023 mean-field gaussian likelihood on the emergent property statistics. The mean and variances of the  
1024 mean-field gaussian likelihood are predicated by  $\boldsymbol{\eta}_{\text{opt}}$  (Equations 35 and 37), which is adapted after  
1025 each EPI optimization epoch based on  $\mathcal{X}$  (see Section 5.1.4). In EPI, the inferred distribution is  
1026 not conditioned on a finite dataset as in variational inference, but rather the emergent property  
1027  $\mathcal{X}$  dictates the likelihood parameterization such that the inferred distribution will produce the  
1028 emergent property. As a note, we could not simply choose  $\boldsymbol{\mu}_f$  and  $\boldsymbol{\sigma}_f$  directly from the outset, since  
1029 we do not know which of these choices will produce the emergent property  $\mathcal{X}$ , which necessitates  
1030 the EPI optimization routine that adapts  $\boldsymbol{\eta}_{\text{opt}}$ . Accordingly, we replace the notation of  $p(\mathbf{z} | \mathbf{x})$   
1031 with  $p(\mathbf{z} | \mathcal{X})$  conceptualizing an inferred distribution that obeys emergent property  $\mathcal{X}$  (see Section  
1032 5.1).  
1033

## 5.2 Stomatogastric ganglion

1034 In Section 3.1 and 3.2, we used EPI to infer conductance parameters in a model of the stomatogastric  
1035 ganglion (STG) [41]. This 5-neuron circuit model represents two subcircuits: that generating the  
1036 pyloric rhythm (fast population) and that generating the gastric mill rhythm (slow population).  
1037 The additional neuron (the IC neuron of the STG) receives inhibitory synaptic input from both  
1038 subcircuits, and can couple to either rhythm dependent on modulatory conditions. There is also  
1039 a parametric regime in which this neuron fires at an intermediate frequency between that of the  
1040 fast and slow populations [41], which we infer with EPI as a motivational example. This model

1041 is not to be confused with an STG subcircuit model of the pyloric rhythm [57], which has been  
 1042 statistically inferred in other studies [15, 35].

1043 **5.2.1 STG model**

1044 We analyze how the parameters  $\mathbf{z} = [g_{el}, g_{synA}]$  govern the emergent phenomena of intermediate  
 1045 hub frequency in a model of the stomatogastric ganglion (STG) [41] shown in Figure 1A with  
 1046 activity  $\mathbf{x} = [x_{f1}, x_{f2}, x_{hub}, x_{s1}, x_{s2}]$ , using the same hyperparameter choices as Gutierrez et al.  
 1047 Each neuron's membrane potential  $x_\alpha(t)$  for  $\alpha \in \{f1, f2, \text{hub}, s1, s2\}$  is the solution of the following  
 1048 stochastic differential equation:

$$C_m \frac{dx_\alpha}{dt} = -[h_{leak}(\mathbf{x}; \mathbf{z}) + h_{Ca}(\mathbf{x}; \mathbf{z}) + h_K(\mathbf{x}; \mathbf{z}) + h_{hyp}(\mathbf{x}; \mathbf{z}) + h_{elec}(\mathbf{x}; \mathbf{z}) + h_{syn}(\mathbf{x}; \mathbf{z})] + dB. \quad (39)$$

1049 The input current of each neuron is the sum of the leak, calcium, potassium, hyperpolarization,  
 1050 electrical and synaptic currents. Each current component is a function of all membrane potentials  
 1051 and the conductance parameters  $\mathbf{z}$ . Finally, we include gaussian noise  $dB$  to the model of Gutierrez  
 1052 et al. so that the model stochastic, although this is not required by EPI.

1053 The capacitance of the cell membrane was set to  $C_m = 1nF$ . Specifically, the currents are the  
 1054 difference in the neuron's membrane potential and that current type's reversal potential multiplied  
 1055 by a conductance:

$$h_{leak}(\mathbf{x}; \mathbf{z}) = g_{leak}(x_\alpha - V_{leak}) \quad (40)$$

$$h_{elec}(\mathbf{x}; \mathbf{z}) = g_{el}(x_\alpha^{post} - x_\alpha^{pre}) \quad (41)$$

$$h_{syn}(\mathbf{x}; \mathbf{z}) = g_{syn}S_\infty^{pre}(x_\alpha^{post} - V_{syn}) \quad (42)$$

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

$$h_K(\mathbf{x}; \mathbf{z}) = g_KN(x_\alpha - V_K) \quad (44)$$

$$h_{hyp}(\mathbf{x}; \mathbf{z}) = g_hH(x_\alpha - V_{hyp}). \quad (45)$$

1061 The reversal potentials were set to  $V_{leak} = -40mV$ ,  $V_{Ca} = 100mV$ ,  $V_K = -80mV$ ,  $V_{hyp} = -20mV$ ,  
 1062 and  $V_{syn} = -75mV$ . The other conductance parameters were fixed to  $g_{leak} = 1 \times 10^{-4}\mu S$ .  $g_{Ca}$ ,  
 1063  $g_K$ , and  $g_{hyp}$  had different values based on fast, intermediate (hub) or slow neuron. The fast  
 1064 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  
 1065 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  
 1066 slow conductances had values  $g_{Ca} = 8.5 \times 10^{-3}$ ,  $g_K = 1.5 \times 10^{-2}$ , and  $g_{hyp} = 1.0 \times 10^{-2}$ .

1067 Furthermore, the Calcium, Potassium, and hyperpolarization channels have time-dependent gating  
 1068 dynamics dependent on steady-state gating variables  $M_\infty$ ,  $N_\infty$  and  $H_\infty$ , respectively:

$$M_\infty = 0.5 \left( 1 + \tanh \left( \frac{x_\alpha - v_1}{v_2} \right) \right) \quad (46)$$

$$\frac{dN}{dt} = \lambda_N (N_\infty - N) \quad (47)$$

$$N_\infty = 0.5 \left( 1 + \tanh \left( \frac{x_\alpha - v_3}{v_4} \right) \right) \quad (48)$$

$$\lambda_N = \phi_N \cosh \left( \frac{x_\alpha - v_3}{2v_4} \right) \quad (49)$$

$$\frac{dH}{dt} = \frac{(H_\infty - H)}{\tau_h} \quad (50)$$

$$H_\infty = \frac{1}{1 + \exp \left( \frac{x_\alpha + v_5}{v_6} \right)} \quad (51)$$

$$\tau_h = 272 - \left( \frac{-1499}{1 + \exp \left( \frac{-x_\alpha + v_7}{v_8} \right)} \right). \quad (52)$$

1075 where we set  $v_1 = 0mV$ ,  $v_2 = 20mV$ ,  $v_3 = 0mV$ ,  $v_4 = 15mV$ ,  $v_5 = 78.3mV$ ,  $v_6 = 10.5mV$ ,  
 1076  $v_7 = -42.2mV$ ,  $v_8 = 87.3mV$ ,  $v_9 = 5mV$ , and  $v_{th} = -25mV$ .

1077 Finally, there is a synaptic gating variable as well:

$$S_\infty = \frac{1}{1 + \exp \left( \frac{v_{th} - x_\alpha}{v_9} \right)}. \quad (53)$$

1078 When the dynamic gating variables are considered, this is actually a 15-dimensional nonlinear  
 1079 dynamical system. The gaussian noise  $d\mathbf{B}$  has variance  $(1 \times 10^{-12})^2$  A<sup>2</sup>, and introduces variability  
 1080 in frequency at each parameterization  $\mathbf{z}$ .

### 1081 5.2.2 Hub frequency calculation

1082 In order to measure the frequency of the hub neuron during EPI, the STG model was simulated for  
 1083  $T = 300$  time steps of  $dt = 25\text{ms}$ . The chosen  $dt$  and  $T$  were the most computationally convenient  
 1084 choices yielding accurate frequency measurement. We used a basis of complex exponentials with  
 1085 frequencies from 0.0-1.0 Hz at 0.01Hz resolution to measure frequency from simulated time series

$$\Phi = [0.0, 0.01, \dots, 1.0]^\top .. \quad (54)$$

1086 To measure spiking frequency, we processed simulated membrane potentials with a relu (spike  
 1087 extraction) and low-pass filter with averaging window of size 20, then took the frequency with the

1088 maximum absolute value of the complex exponential basis coefficients of the processed time-series.  
 1089 The first 20 temporal samples of the simulation are ignored to account for initial transients.  
 1090 To differentiate through the maximum frequency identification, we used a soft-argmax Let  $X_\alpha \in$   
 1091  $\mathcal{C}^{|\Phi|}$  be the complex exponential filter bank dot products with the signal  $x_\alpha \in \mathbb{R}^N$ , where  $\alpha \in$   
 1092  $\{f1, f2, \text{hub}, s1, s2\}$ . The soft-argmax is then calculated using temperature parameter  $\beta = 100$

$$\psi_\alpha = \text{softmax}(\beta |X_\alpha| \odot i), \quad (55)$$

1093 where  $i = [0, 1, \dots, 100]$ . The frequency is then calculated as

$$\omega_\alpha = 0.01\psi_\alpha \text{Hz}. \quad (56)$$

1094 Intermediate hub frequency, like all other emergent properties in this work, is defined by the mean  
 1095 and variance of the emergent property statistics. In this case, we have one statistic, hub neuron  
 1096 frequency, where the mean was chosen to be 0.55Hz,(Equation 2) and variance was chosen to be  
 1097  $(0.025\text{Hz})^2$  (Equation 3).

### 1098 5.2.3 EPI details for the STG model

1099 As a maximum entropy distribution,  $T(\mathbf{x}; \mathbf{z})$  is comprised of both these first and second moments  
 1100 of the hub neuron frequency (as in Equations 17 and 18)

$$T(\mathbf{x}; \mathbf{z}) = \begin{bmatrix} \omega_{\text{hub}}(\mathbf{x}; \mathbf{z}) \\ (\omega_{\text{hub}}(\mathbf{x}; \mathbf{z}) - 0.55)^2 \end{bmatrix}, \quad (57)$$

$$\boldsymbol{\mu}_{\text{opt}} = \begin{bmatrix} 0.55 \\ 0.025^2 \end{bmatrix}. \quad (58)$$

1101  
 1102 Throughout optimization, the augmented Lagrangian parameters  $\eta$  and  $c$ , were updated after each  
 1103 epoch of 5,000 iterations(see Section 5.1.4). The optimization converged after five epochs (Fig. S4).  
 1104 For EPI in Fig 1E, we used a real NVP architecture with three Real NVP coupling layers and two-  
 1105 layer neural networks of 25 units per layer. The normalizing flow architecture mapped  $\mathbf{z}_0 \sim \mathcal{N}(\mathbf{0}, I)$   
 1106 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  
 1107 returned by a preliminary ABC search. We did not include  $g_{\text{synA}} < 0.01$ , for numerical stability.  
 1108 EPI optimization was run using 5 different random seeds for architecture initialization  $\boldsymbol{\theta}$  with an  
 1109 augmented Lagrangian coefficient of  $c_0 = 10^5$ , a batch size  $n = 400$ , and  $\beta = 2$ . The architecture  
 1110 converged with criteria  $N_{\text{test}} = 100$ .

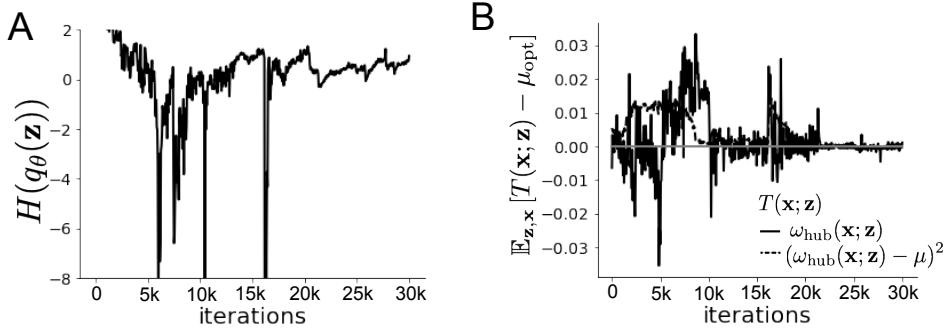


Figure 9: (STG1): EPI optimization of the STG model producing network syncing. **A.** Entropy throughout optimization. **B.** The emergent property statistic means and variances converge to their constraints at 25,000 iterations following the fifth augmented Lagrangian epoch.

#### 1111 5.2.4 Hessian sensitivity vectors

1112 To quantify the second-order structure of the EPI distribution, we evaluated the Hessian of the log  
 1113 probability  $\frac{\partial^2 \log q(\mathbf{z}|\mathcal{X})}{\partial \mathbf{z} \mathbf{z}^\top}$ . The eigenvector of this Hessian with most negative eigenvalue is defined as  
 1114 the sensitivity dimension  $\mathbf{v}_1$ , and all subsequent eigenvectors are ordered by increasing eigenvalue.  
 1115 These eigenvalues are quantifications of how fast the emergent property deteriorates via the param-  
 1116 eter combination of their associated eigenvector. In Figure 1D, the sensitivity dimension  $v_1$  (solid)  
 1117 and the second eigenvector of the Hessian  $v_2$  (dashed) are shown evaluated at the mode of the  
 1118 distribution. Since the Hessian eigenvectors have sign degeneracy, the visualized directions in 2-D  
 1119 parameter space were chosen to have positive  $g_{\text{synA}}$ . The length of the arrows is inversely propor-  
 1120 tional to the square root of the absolute value of their eigenvalues  $\lambda_1 = -10.7$  and  $\lambda_2 = -3.22$ . For  
 1121 the same magnitude perturbation away from the mode, intermediate hub frequency only diminishes  
 1122 along the sensitivity dimension  $\mathbf{v}_1$  (Fig. 1E-F).

#### 1123 5.3 Scaling EPI for stable amplification in RNNs

##### 1124 5.3.1 Rank-2 RNN model

1125 We examined the scaling properties of EPI by learning connectivities of RNNs of increasing size  
 1126 that exhibit stable amplification. Rank-2 RNN connectivity was modeled as  $W = UV^\top$ , where  
 1127  $U = [\mathbf{U}_1 \ \mathbf{U}_2] + g\chi^{(W)}$ ,  $V = [\mathbf{V}_1 \ \mathbf{V}_2] + g\chi^{(V)}$ , and  $\chi_{i,j}^{(W)}, \chi_{i,j}^{(V)} \sim \mathcal{N}(0, 1)$ . This RNN model has

<sub>1128</sub> dynamics

$$\tau \dot{\mathbf{x}} = -\mathbf{x} + W\mathbf{x}. \quad (59)$$

<sub>1129</sub> In this analysis, we inferred connectivity parameterizations  $\mathbf{z} = [\mathbf{U}_1^\top, \mathbf{U}_2^\top, \mathbf{V}_1^\top, \mathbf{V}_2^\top]^\top \in [-1, 1]^{(4N)}$   
<sub>1130</sub> that produced stable amplification using EPI, SMC-ABC [26], and SNPE [35] (see Section Related  
<sub>1131</sub> Methods).

### <sub>1132</sub> 5.3.2 Stable amplification

<sub>1133</sub> For this RNN model to be stable, all real eigenvalues of  $W$  must be less than 1:  $\text{real}(\lambda_1) < 1$ ,  
<sub>1134</sub> where  $\lambda_1$  denotes the greatest real eigenvalue of  $W$ . For a stable RNN to amplify at least one input  
<sub>1135</sub> pattern, the symmetric connectivity  $W^s = \frac{W+W^\top}{2}$  must have an eigenvalue greater than 1:  $\lambda_1^s > 1$ ,  
<sub>1136</sub> where  $\lambda^s$  is the maximum eigenvalue of  $W^s$ . These two conditions are necessary and sufficient for  
<sub>1137</sub> stable amplification in RNNs [51].

### <sub>1138</sub> 5.3.3 EPI details for RNNs

<sub>1139</sub> We defined the emergent property of stable amplification with means of these eigenvalues (0.5  
<sub>1140</sub> and 1.5, respectively) that satisfy these conditions. To complete the emergent property definition,  
<sub>1141</sub> we chose variances ( $0.25^2$ ) about those means such that samples rarely violate the eigenvalue  
<sub>1142</sub> constraints. In terms of the EPI optimization variables, this is written as

$$T(\mathbf{x}; \mathbf{z}) = \begin{bmatrix} \text{real}(\lambda_1)(\mathbf{x}; \mathbf{z}) \\ \lambda_1^s(\mathbf{x}; \mathbf{z}) \\ (\text{real}(\lambda_1)(\mathbf{x}; \mathbf{z}) - 0.5)^2 \\ (\lambda_1^s(\mathbf{x}; \mathbf{z}) - 1.5)^2 \end{bmatrix}, \quad (60)$$

<sub>1143</sub>

$$\boldsymbol{\mu}_{\text{opt}} = \begin{bmatrix} 0.5 \\ 1.5 \\ 0.25^2 \\ 0.25^2 \end{bmatrix}. \quad (61)$$

<sub>1144</sub> Gradients of maximum eigenvalues of Hermitian matrices like  $W^s$  are available with modern auto-  
<sub>1145</sub> matic differentiation tools. To differentiate through the  $\text{real}(\lambda_1)$ , we solved the following equation  
<sub>1146</sub> for eigenvalues of rank-2 matrices using the rank reduced matrix  $W^r = V^\top U$

$$\lambda_{\pm} = \frac{\text{Tr}(W^r) \pm \sqrt{\text{Tr}(W^r)^2 - 4\text{Det}(W^r)}}{2}. \quad (62)$$

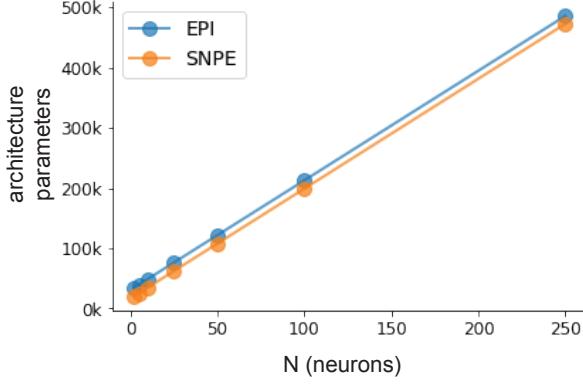


Figure 10: (RNN1): Number of parameters in deep probability distribution architectures of EPI (blue) and SNPE (orange) by RNN size ( $N$ ).

1147 For EPI in Fig. 2, we used a real NVP architecture with three coupling layers of affine transformations parameterized by two-layer neural networks of 100 units per layer. The initial distribution  
 1148 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  
 1149 an augmented Lagrangian coefficient of  $c_0 = 10^3$ , a batch size  $n = 200$ ,  $\beta = 4$ , and chose to use  
 1150 500 iterations per augmented Lagrangian epoch and emergent property constraint convergence was  
 1151 evaluated at  $N_{\text{test}} = 200$  (Fig. 2B blue line, and Fig. 2C-D blue).

1153 **5.3.4 Methodological comparison**

1154 We compared EPI to two alternative simulation-based inference techniques, since the likelihood  
 1155 of these eigenvalues given  $\mathbf{z}$  is not available. Approximate Bayesian computation (ABC) [24] is a  
 1156 rejection sampling technique for obtaining sets of parameters  $\mathbf{z}$  that produce activity  $\mathbf{x}$  close to some  
 1157 observed data  $\mathbf{x}_0$ . Sequential Monte Carlo approximate Bayesian computation (SMC-ABC) is the  
 1158 state-of-the-art ABC method, which leverages SMC techniques to improve sampling speed. We ran  
 1159 SMC-ABC with the pyABC package [94] to infer RNNs with stable amplification: connectivities  
 1160 having eigenvalues within an  $\epsilon$ -defined  $l_2$  distance of

$$x_0 = \begin{bmatrix} \text{real}(\lambda_1) \\ \lambda_1^s \end{bmatrix} = \begin{bmatrix} 0.5 \\ 1.5 \end{bmatrix}. \quad (63)$$

1161 SMC-ABC was run with a uniform prior over  $\mathbf{z} \in [-1, 1]^{(4N)}$ , a population size of 1,000 particles  
 1162 with simulations parallelized over 32 cores, and a multivariate normal transition model.

1163 SNPE, the next approach in our comparison, is far more similar to EPI. Like EPI, SNPE treats pa-

1164 rameters in mechanistic models with deep probability distributions, yet the two learning algorithms  
 1165 are categorically different. SNPE uses a two-network architecture to approximate the posterior dis-  
 1166 tribution of the model conditioned on observed data  $\mathbf{x}_0$ . The amortizing network maps observations  
 1167  $\mathbf{x}_i$  to the parameters of the deep probability distribution. The weights and biases of the parameter  
 1168 network are optimized by sequentially augmenting the training data with additional pairs  $(\mathbf{z}_i, \mathbf{x}_i)$   
 1169 based on the most recent posterior approximation. This sequential procedure is important to get  
 1170 training data  $\mathbf{z}_i$  to be closer to the true posterior, and  $\mathbf{x}_i$  to be closer to the observed data. For  
 1171 the deep probability distribution architecture, we chose a masked autoregressive flow with affine  
 1172 couplings (the default choice), three transforms, 50 hidden units, and a normalizing flow mapping  
 1173 to the support as in EPI. This architectural choice closely tracked the size of the architecture used  
 1174 by EPI (Fig. 10). As in SMC-ABC, we ran SNPE with  $\mathbf{x}_0 = \mu$ . All SNPE optimizations were  
 1175 run for a limit of 1.5 days on a Tesla V100 GPU, or until two consecutive rounds resulted in a  
 1176 validation log probability lower than the maximum observed for that random seed.

1177 To compare the efficiency of these algorithms for inferring RNN connectivity distributions producing  
 1178 stable amplification, we develop a convergence criteria that can be used across methods. While EPI  
 1179 has its own hypothesis testing convergence criteria for the emergent property, it would not make  
 1180 sense to use this criteria on SNPE and SMC-ABC which do not constrain the means and variances  
 1181 of their predictions. Instead, we consider EPI and SNPE to have converged after completing its  
 1182 most recent optimization epoch (EPI) or round (SNPE) in which the distance

$$d(q_\theta(z)) = |\mathbb{E}_{\mathbf{z}, \mathbf{x}} [f(\mathbf{x}; \mathbf{z})] - \boldsymbol{\mu}|_2 \quad (64)$$

1183 is less than 0.5. We consider SMC-ABC to have converged once the population produces samples  
 1184 within the  $\epsilon = 0.5$  ball ensuring stable amplification.

1185 When assessing the scalability of SNPE, it is important to check that alternative hyperparamter-  
 1186 izations could not yield better performance. Key hyperparameters of the SNPE optimization are  
 1187 the number of simulations per round  $n_{\text{round}}$ , the number of atoms used in the atomic proposals of  
 1188 the SNPE-C algorithm [95], and the batch size  $n$ . To match EPI, we used a batch size of  $n = 200$   
 1189 for  $N <= 25$ , however we found  $n = 1,000$  to be helpful for SNPE in higher dimensions. While  
 1190  $n_{\text{round}} = 1,000$  yielded SNPE convergence for  $N <= 25$ , we found that a substantial increase to  
 1191  $n_{\text{round}} = 25,000$  yielded more consistent convergence at  $N = 50$  (Fig. 11A). By increasing  $n_{\text{round}}$ ,  
 1192 we also necessarily increase the duration of each round. At  $N = 100$ , we tried two hyperparameter  
 1193 modifications. As suggested in [95], we increased  $n_{\text{atom}}$  by an order of magnitude to improve gra-  
 1194 dient quality, but this had little effect on the optimization (much overlap between same random

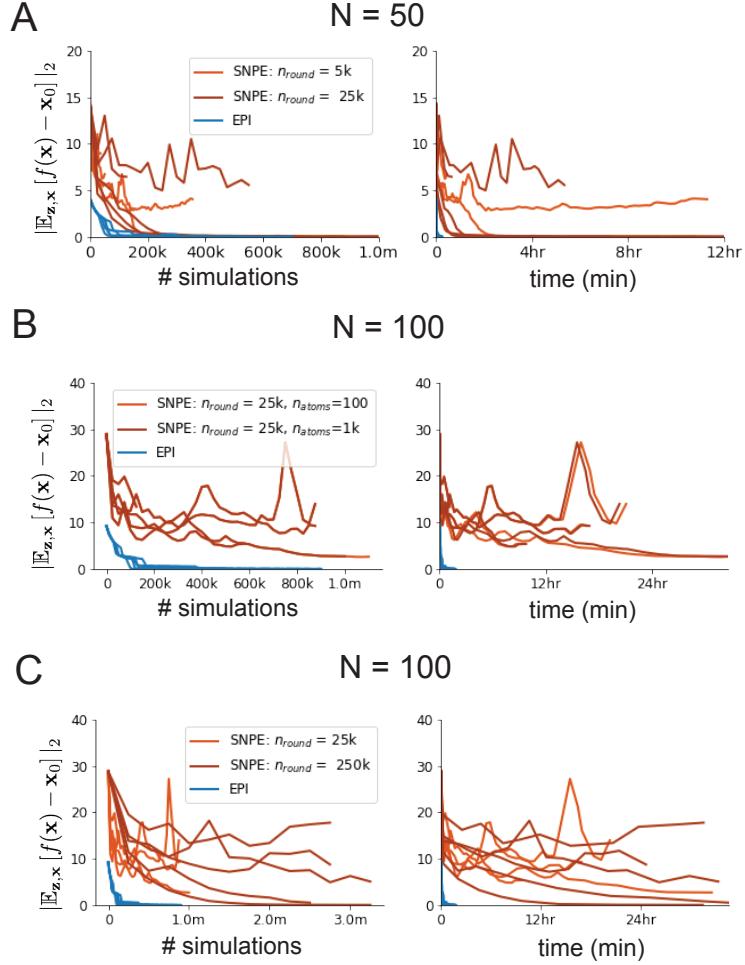


Figure 11: (RNN3): SNPE convergence was enabled by increasing  $n_{\text{round}}$ , not  $n_{\text{atom}}$ . **A.** Difference of mean predictions  $\mathbf{x}_0$  throughout optimization at  $N = 50$  with by simulation count (left) and wall time (right) of SNPE with  $n_{\text{round}} = 5,000$  (light 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).

seeds) (Fig. 11B). Finally, we increased  $n_{\text{round}}$  by an order of magnitude, which yielded convergence in one case, but no others. We found no way to improve the convergence rate of SNPE without making more aggressive hyperparameter choices requiring high numbers of simulations. In Figure 2C-D, we show samples from the random seed resulting in emergent property convergence at greatest entropy (EPI), the random seed resulting in greatest validation log probability (SNPE), and the result of all converged random seeds (SMC).

### 5.3.5 Effect of RNN parameters on EPI and SNPE inferred distributions

To clarify the difference in objectives of EPI and SNPE, we show their results on RNN models with different numbers of neurons  $N$  and random strength  $g$ . The parameters inferred by EPI consistently produces the same mean and variance of  $\text{real}(\lambda_1)$  and  $\lambda_1^s$ , while those inferred by SNPE change according to the model definition (Fig. 12A). For  $N = 2$  and  $g = 0.01$ , the SNPE posterior has greater concentration in eigenvalues around  $\mathbf{x}_0$  than at  $g = 0.1$ , where the model has greater randomness (Fig. 12B top, orange). At both levels of  $g$  when  $N = 2$ , the posterior of SNPE has lower entropy than EPI at convergence (Fig. 12B top). However at  $N = 10$ , SNPE results in a predictive distribution of more widely dispersed eigenvalues (Fig. 12A bottom), and an inferred posterior with greater entropy than EPI (Fig. 12B bottom). We highlight these differences not to focus on an insightful trend, but to emphasize that these methods optimize different objectives with different implications.

Note that SNPE converges when it's validation log probability has saturated after several rounds of optimization (Fig. 12C), and that EPI converges after several epochs of its own optimization to enforce the emergent property constraints (Fig. 12D blue). Importantly, as SNPE optimizes its posterior approximation, the predictive means change, and at convergence may be different than  $\mathbf{x}_0$  (Fig. 12D orange, left). It is sensible to assume that predictions of a well-approximated SNPE posterior should closely reflect the data on average (especially given a uniform prior and a low degree of stochasticity), however this is not a given. Furthermore, no aspect of the SNPE optimization controls the variance of the predictions (Fig. 12D orange, right).

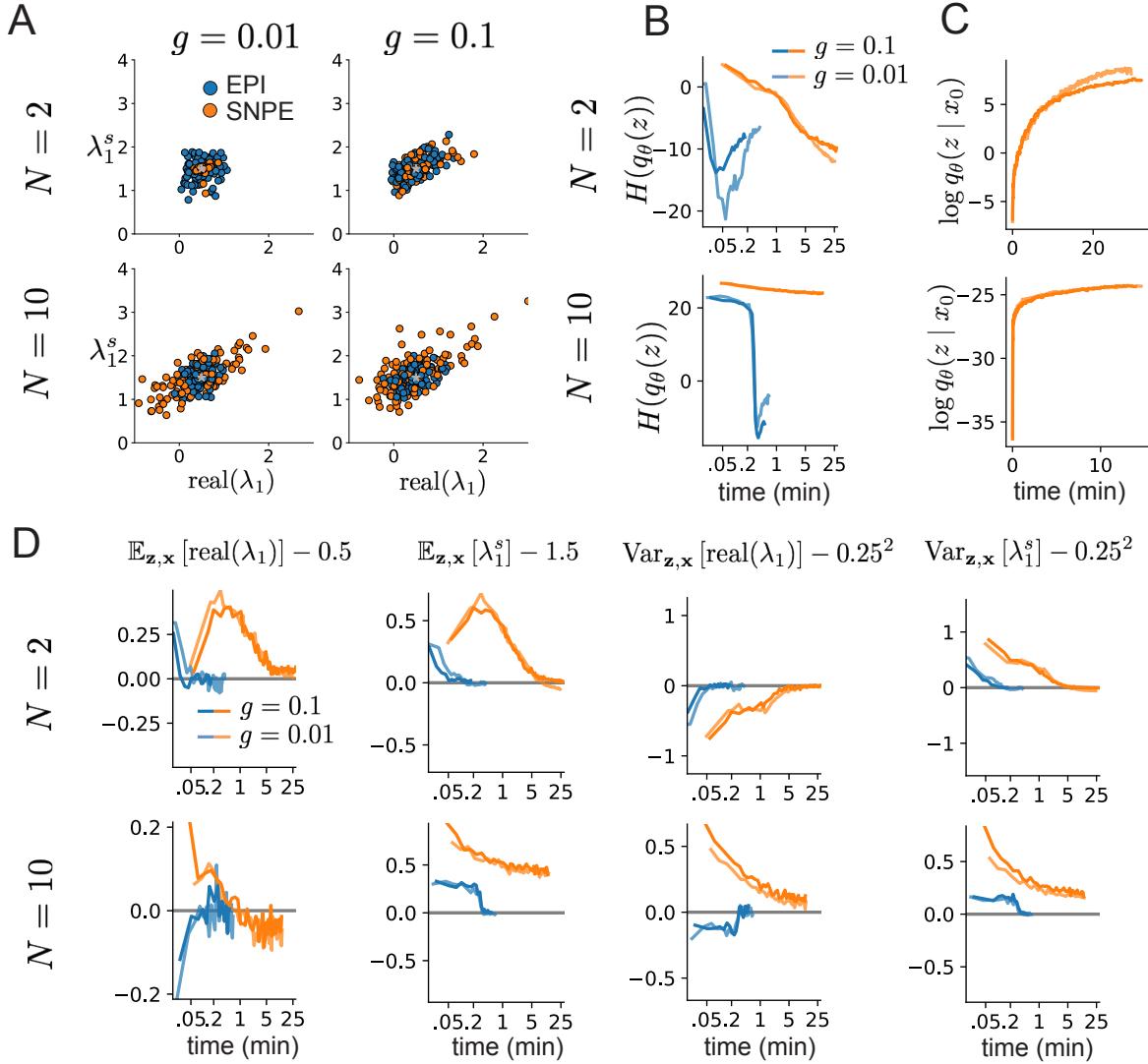


Figure 12: (RNN2): Model characteristics affect predictions of posteriors inferred by SNPE, while predictions of parameters inferred by EPI remain fixed. **A.** Predictive distribution of EPI (blue) and SNPE (orange) inferred connectivity of RNNs exhibiting stable amplification with  $N = 2$  (top),  $N = 10$  (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.

1221 **5.4 Primary visual cortex**

1222 **5.4.1 V1 model**

1223 E-I circuit models, rely on the assumption that inhibition can be studied as an indivisible unit,  
 1224 despite ample experimental evidence showing that inhibition is instead composed of distinct ele-  
 1225 ments [65]. In particular three types of genetically identified inhibitory cell-types – parvalbumin  
 1226 (P), somatostatin (S), VIP (V) – compose 80% of GABAergic interneurons in V1 [63–65], and follow  
 1227 specific connectivity patterns (Fig. 3A) [66], which lead to cell-type specific computations [47, 96].  
 1228 Currently, how the subdivision of inhibitory cell-types, shapes correlated variability by reconfigur-  
 1229 ing recurrent network dynamics is not understood.

1230 In the stochastic stabilized supralinear network [60], population rate responses  $\mathbf{x}$  to mean input  $\mathbf{h}$ ,  
 1231 recurrent input  $W\mathbf{x}$  and slow noise  $\epsilon$  are governed by

$$\tau \frac{d\mathbf{x}}{dt} = -\mathbf{x} + \phi(W\mathbf{x} + \mathbf{h} + \epsilon), \quad (65)$$

1232 where the noise is an Ornstein-Uhlenbeck process  $\epsilon \sim OU(\tau_{\text{noise}}, \sigma)$

$$\tau_{\text{noise}} d\epsilon_\alpha = -\epsilon_\alpha dt + \sqrt{2\tau_{\text{noise}}} \tilde{\sigma}_\alpha dB \quad (66)$$

1233 with  $\tau_{\text{noise}} = 5\text{ms} > \tau = 1\text{ms}$ . The noisy process is parameterized as

$$\tilde{\sigma}_\alpha = \sigma_\alpha \sqrt{1 + \frac{\tau}{\tau_{\text{noise}}}}, \quad (67)$$

1234 so that  $\sigma$  parameterizes the variance of the noisy input in the absence of recurrent connectivity  
 1235 ( $W = \mathbf{0}$ ). As contrast  $c \in [0, 1]$  increases, input to the E- and P-populations increases relative to  
 1236 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  
 1237 using the deterministic V1 circuit model [47]

$$W_{\text{fit}} = \begin{bmatrix} W_{EE} & W_{EP} & W_{ES} & W_{EV} \\ W_{PE} & W_{PP} & W_{PS} & W_{PV} \\ W_{SE} & W_{SP} & W_{SS} & W_{SV} \\ W_{VE} & W_{VP} & W_{VS} & W_{VV} \end{bmatrix} = \begin{bmatrix} 2.18 & -1.19 & -.594 & -.229 \\ 1.66 & -.651 & -.680 & -.242 \\ .895 & -5.22 \times 10^{-3} & -1.51 \times 10^{-4} & -.761 \\ 3.34 & -2.31 & -.254 & -2.52 \times 10^{-4} \end{bmatrix}, \quad (68)$$

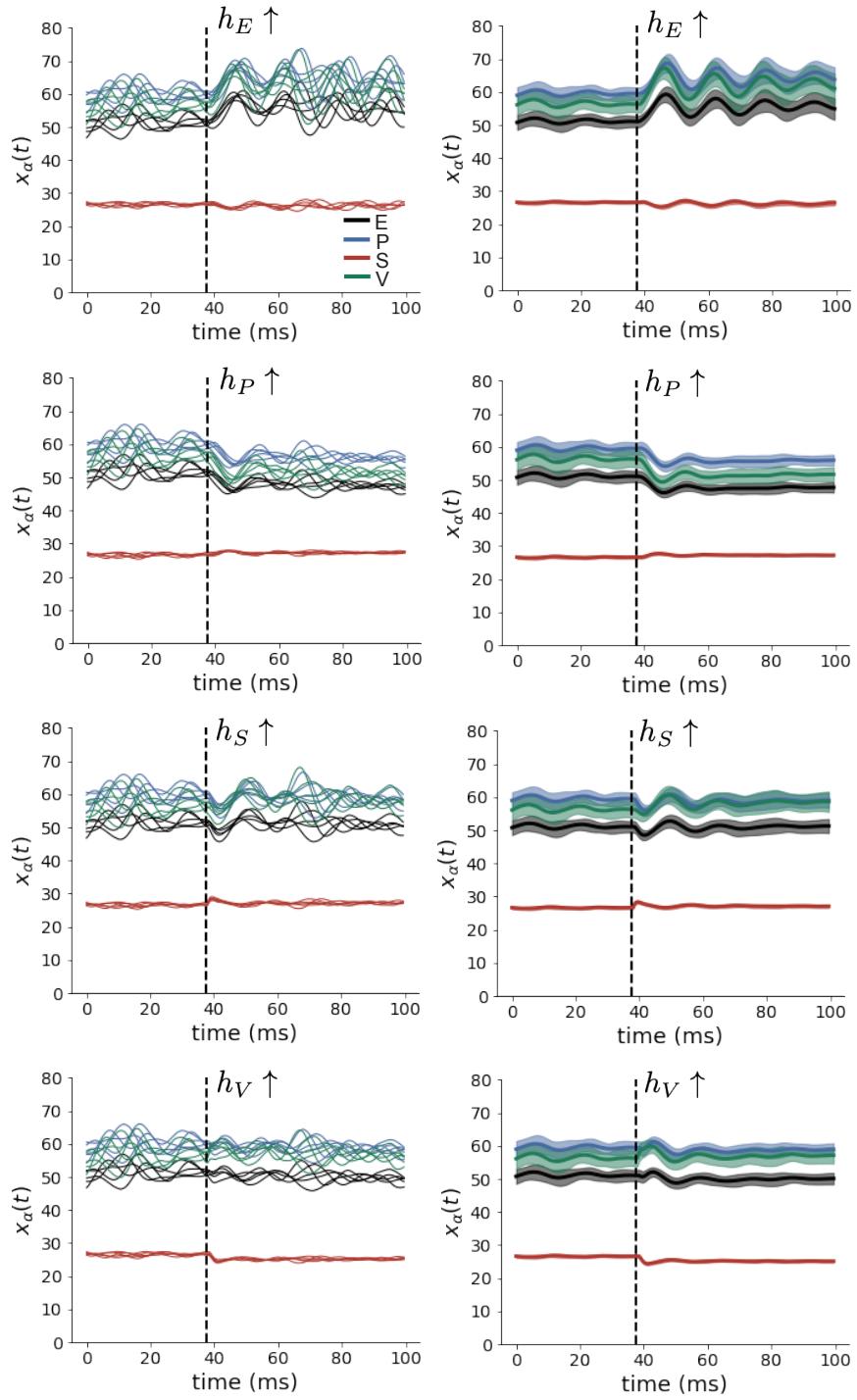


Figure 13: (V1 1) (Left) Simulations for small increases in neuron-type population input. Input magnitudes are chosen so that effect is salient (0.002 for E and P, but 0.02 for S and V). (Right) Average (solid) and standard deviation (shaded) of stochastic fluctuations of responses.

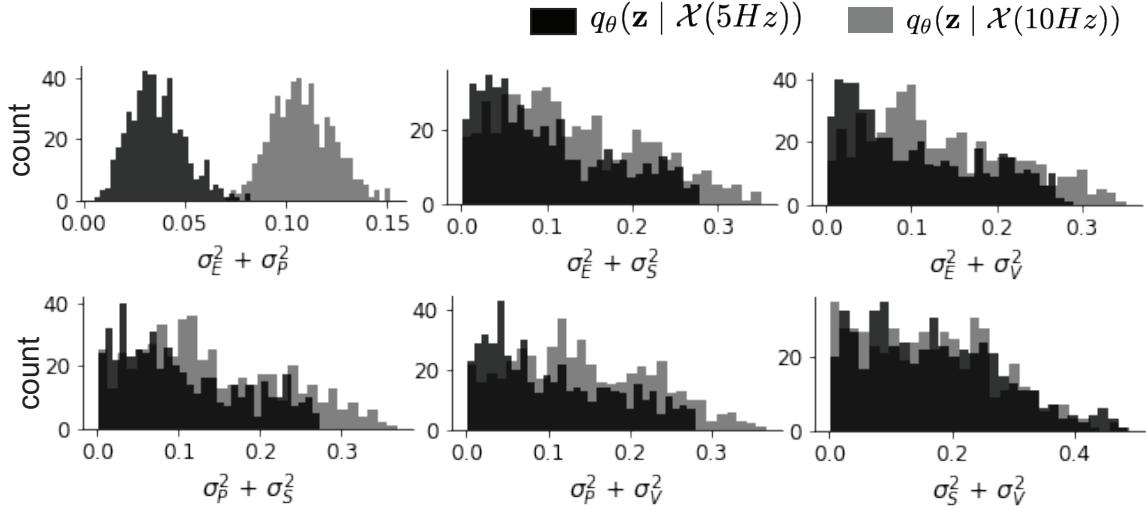


Figure 14: (V1 2) EPI predictive distributions of the sum of squares of each pair of noise parameters.

$$\mathbf{h}_{b,\text{fit}} = \begin{bmatrix} .416 \\ .429 \\ .491 \\ .486 \end{bmatrix}, \quad (69)$$

<sup>1238</sup> and

$$\mathbf{h}_{c,\text{fit}} = \begin{bmatrix} .359 \\ .403 \\ 0 \\ 0 \end{bmatrix}. \quad (70)$$

<sup>1239</sup> To obtain rates on a realistic scale (100-fold greater), we map these fitted parameters to an equivalence class  
<sup>1240</sup>

$$W = \begin{bmatrix} W_{EE} & W_{EP} & W_{ES} & W_{EV} \\ W_{PE} & W_{PP} & W_{PS} & W_{PV} \\ W_{SE} & W_{SP} & W_{SS} & W_{SV} \\ W_{VE} & W_{VP} & W_{VS} & W_{VV} \end{bmatrix} = \begin{bmatrix} .218 & -.119 & -.0594 & -.0229 \\ .166 & -.0651 & -.068 & -.0242 \\ .0895 & -5.22 \times 10^{-4} & -1.51 \times 10^{-5} & -.0761 \\ .334 & -.231 & -.0254 & -2.52 \times 10^{-5} \end{bmatrix}, \quad (71)$$

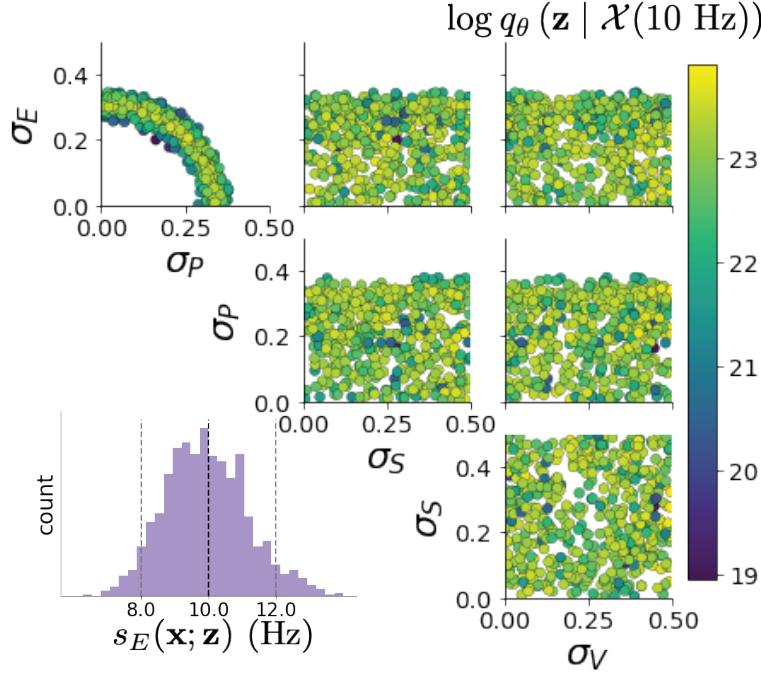


Figure 15: (V1 3) EPI inferred distribution for  $\mathcal{X}(10 \text{ Hz})$ .

$$\mathbf{h}_b = \begin{bmatrix} h_{b,E} \\ h_{b,P} \\ h_{b,S} \\ h_{b,V} \end{bmatrix} = \begin{bmatrix} 4.16 \\ 4.29 \\ 4.91 \\ 4.86 \end{bmatrix}, \quad (72)$$

<sup>1241</sup> and

$$\mathbf{h}_c = \begin{bmatrix} h_{c,E} \\ h_{c,P} \\ h_{c,S} \\ h_{c,V} \end{bmatrix} = \begin{bmatrix} 3.59 \\ 4.03 \\ 0 \\ 0 \end{bmatrix}. \quad (73)$$

<sup>1242</sup> Circuit responses are simulated using  $T = 200$  time steps at  $dt = 0.5\text{ms}$  from an initial condition  
<sup>1243</sup> 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  
<sup>1244</sup> as the square root of the temporal variance from  $t_{ss} = 75\text{ms}$  to  $Tdt = 100\text{ms}$  averaged over 100  
<sup>1245</sup> independent trials.

$$s_E(\mathbf{x}; \mathbf{z}) = \mathbb{E}_x \left[ \sqrt{\mathbb{E}_{t > t_{ss}} [(x_E(t) - \mathbb{E}_{t > t_{ss}} [x_E(t)])^2]} \right] \quad (74)$$

1246 **5.4.2 EPI details for the V1 model**

1247 For EPI in Fig 3D-E, we used a real NVP architecture with three Real NVP coupling layers  
 1248 and two-layer neural networks of 50 units per layer. The normalizing flow architecture mapped  
 1249  $z_0 \sim \mathcal{N}(\mathbf{0}, I)$  to a support of  $\mathbf{z} = [\sigma_E, \sigma_P, \sigma_S, \sigma_V] \in [0.0, 0.5]^4$ . EPI optimization was run using three  
 1250 different random seeds for architecture initialization  $\boldsymbol{\theta}$  with an augmented Lagrangian coefficient of  
 1251  $c_0 = 10^{-1}$ , a batch size  $n = 100$ , and  $\beta = 2$ . The distributions shown are those of the architectures  
 1252 converging with criteria  $N_{\text{test}} = 100$  at greatest entropy across three random seeds.

1253 **5.4.3 Sensitivity analyses**

1254 In Fig. 3E, we visualize the modes of  $q_{\boldsymbol{\theta}}(\mathbf{z} \mid \mathcal{X})$  throughout the  $\sigma_E$ - $\sigma_P$  marginal. Specifically, we  
 1255 calculated

$$\begin{aligned} \mathbf{z}^*(\sigma_{P,\text{fixed}}) &= \underset{\mathbf{z}}{\operatorname{argmax}} \log q_{\boldsymbol{\theta}}(\mathbf{z} \mid \mathcal{X}) \\ \text{s.t. } \sigma_P &= \sigma_{P,\text{fixed}} \end{aligned} \quad (75)$$

1256 At each mode  $\mathbf{z}^*$ , we calculated the Hessian and visualized the sensitivity dimension in the direction  
 1257 of positive  $\sigma_E$ .

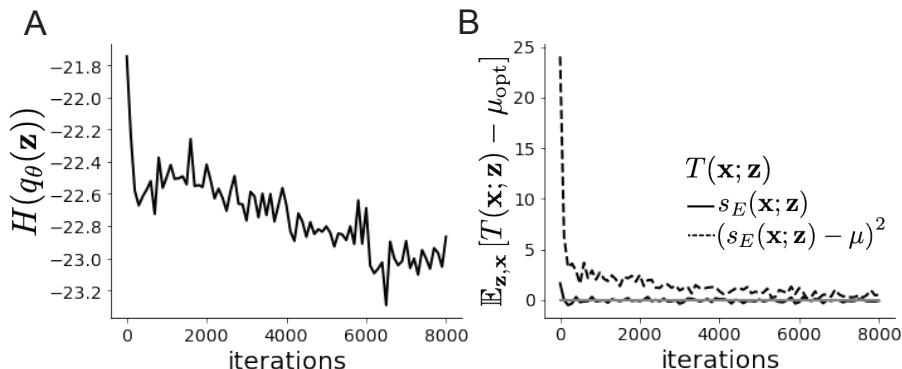


Figure 16: (V1 4) Optimization for V1

1258 **5.4.4 Primary visual cortex: Mathematical intuition and challenges**

1259 The dynamical system that we are working with can be written as

$$\begin{aligned} dx &= \frac{1}{\tau}(-x + f(Wx + h + \epsilon))dt \\ d\epsilon &= -\frac{dt}{\tau_{\text{noise}}} \epsilon + \frac{\sqrt{2}}{\sqrt{\tau_{\text{noise}}}} \Sigma_\epsilon dW \end{aligned} \quad (76)$$

1260 Where in this paper we chose

$$\Sigma_\epsilon = \tau_{\text{noise}} \begin{bmatrix} \tilde{\sigma}_E & 0 & 0 & 0 \\ 0 & \tilde{\sigma}_P & 0 & 0 \\ 0 & 0 & \tilde{\sigma}_S & 0 \\ 0 & 0 & 0 & \tilde{\sigma}_V \end{bmatrix} \quad (77)$$

1261 where  $\tilde{\sigma}_\alpha$  is the reparameterized standard deviation of the noise for population  $\alpha$  from Equation  
1262 67.

1263 In order to compute this covariance, we define  $v = \omega x + h + \epsilon$  and  $S = I - \omega f'(v)$ , to re-write Eq.  
1264 (76) as an 8-dimensional system:

$$d \begin{pmatrix} \delta v \\ \epsilon \end{pmatrix} = - \begin{pmatrix} S & -\frac{\tau_{\text{noise}} - \tau}{\tau \tau_{\text{noise}}} I \\ 0 & \frac{1}{\tau_{\text{noise}}} I \end{pmatrix} \begin{pmatrix} \delta v \\ \epsilon \end{pmatrix} dt + \begin{pmatrix} 0 & \frac{\sqrt{2}}{\sqrt{\tau_{\text{noise}}}} \Sigma_\epsilon \\ 0 & \frac{\sqrt{2}}{\sqrt{\tau_{\text{noise}}}} \Sigma_\epsilon \end{pmatrix} d\mathbf{W} \quad (78)$$

1265 Where  $d\mathbf{W}$  is a vector with the private noise of each variable. The  $d\mathbf{W}$  term is multiplied by a  
1266 non-diagonal matrix is because the noise that the voltage receives is the exact same than the one  
1267 that comes from the OU process and not another process. The solution of this problem is given by  
1268 the Lyapunov Equation [60, 68]:

$$\begin{pmatrix} S & -\frac{\tau_{\text{noise}} - \tau}{\tau \tau_{\text{noise}}} I \\ 0 & \frac{1}{\tau_{\text{noise}}} I \end{pmatrix} \begin{pmatrix} \Lambda_v & \Lambda_c \\ \Lambda_c^T & \Lambda_\epsilon \end{pmatrix} + \begin{pmatrix} \Lambda_v & \Lambda_c \\ \Lambda_c^T & \Lambda_\epsilon \end{pmatrix} \begin{pmatrix} S^T & 0 \\ -\frac{\tau_{\text{noise}} - \tau}{\tau \tau_{\text{noise}}} I & \frac{1}{\tau_{\text{noise}}} I \end{pmatrix} = \begin{pmatrix} \frac{2}{\tau_{\text{noise}}} \Lambda_\epsilon & \frac{2}{\tau_{\text{noise}}} \Lambda_\epsilon \\ \frac{2}{\tau_{\text{noise}}} \Lambda_\epsilon & \frac{2}{\tau_{\text{noise}}} \Lambda_\epsilon \end{pmatrix} \quad (79)$$

1269 To obtain an equation for  $\Lambda_v$ , we solve this block matrix multiplication:

$$S\Lambda_v + \Lambda_v S^T = \frac{2\Lambda_\epsilon}{\tau_{\text{noise}}} + \frac{\tau_{\text{noise}}^2 - \tau^2}{(\tau \tau_{\text{noise}})^2} \left( \left( \frac{1}{\tau_{\text{noise}}} I + S \right)^{-1} \Lambda_\epsilon + \Lambda_\epsilon \left( \frac{1}{\tau_{\text{noise}}} I + S^T \right)^{-1} \right) \quad (80)$$

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 given  $S$ . Due to symmetry arguments [68], in this case the diagonal elements of the covariance matrix of the voltage will have the form:

$$\Lambda_{v_{ii}} = \sum_{i=\{E,P,S,V\}} g_i(S) \sigma_{ii}^2 \quad (81)$$

1270 These coefficients  $g_i(S)$  are complicated functions of the Jacobian of the system. Although expres-  
 1271 sions for these coefficients can be found explicitly, only numerical evaluation of those expressions  
 1272 determine which components of the noisy input are going to strongly influence the variability of ex-  
 1273 citatory population. Showing the generality of this dependence in more complicated noise scenarios  
 1274 (e.g.  $\tau_{\text{noise}} > \tau$  as in Section 3.4), is the focus of current research.

## 1275 5.5 Superior colliculus

### 1276 5.5.1 SC model

1277 The ability to switch between two separate tasks throughout randomly interleaved trials, or “rapid  
 1278 task switching,” has been studied in rats, and midbrain superior colliculus (SC) has been show to  
 1279 play an important in this computation [69]. Neural recordings in SC exhibited two populations of  
 1280 neurons that simultaneously represented both task context (Pro or Anti) and motor response (con-  
 1281 tralateral or ipsilateral to the recorded side), which led to the distinction of two functional classes:  
 1282 the Pro/Contra and Anti/Ipsi neurons [48]. Given this evidence, Duan et al. proposed a model  
 1283 with four functionally-defined neuron-type populations: two in each hemisphere corresponding to  
 1284 the Pro/Contra and Anti/Ipsi populations. We study how the connectivity of this neural circuit  
 1285 governs rapid task switching ability.

1286 The four populations of this model are denoted as left Pro (LP), left Anti (LA), right Pro (RP)  
 1287 and right Anti (RA). Each unit has an activity ( $x_\alpha$ ) and internal variable ( $u_\alpha$ ) related by

$$x_\alpha = \phi(u_\alpha) = \left( \frac{1}{2} \tanh \left( \frac{u_\alpha - a}{b} \right) + \frac{1}{2} \right), \quad (82)$$

1288 where  $\alpha \in \{LP, LA, RA, RP\}$ ,  $a = 0.05$  and  $b = 0.5$  control the position and shape of the nonlin-

1289 earity. We order the neural populations of  $x$  and  $u$  in the following manner

$$\mathbf{x} = \begin{bmatrix} x_{LP} \\ x_{LA} \\ x_{RP} \\ x_{RA} \end{bmatrix} \quad \mathbf{u} = \begin{bmatrix} u_{LP} \\ u_{LA} \\ u_{RP} \\ u_{RA} \end{bmatrix}, \quad (83)$$

1290 which evolve according to

$$\tau \frac{d\mathbf{u}}{dt} = -\mathbf{u} + W\mathbf{x} + \mathbf{h} + d\mathbf{B}. \quad (84)$$

1291 with time constant  $\tau = 0.09s$ , step size 24ms and Gaussian noise  $d\mathbf{B}$  of variance  $0.2^2$ . These  
1292 hyperparameter values are motivated by modeling choices and results from [48].

1293 The weight matrix has 4 parameters for self  $sW$ , vertical  $vW$ , horizontal  $hW$ , and diagonal  $dW$   
1294 connections:

$$W = \begin{bmatrix} sW & vW & hW & dW \\ vW & sW & dW & hW \\ hW & dW & sW & vW \\ dW & hW & vW & sW \end{bmatrix}. \quad (85)$$

1295 We study the role of parameters  $\mathbf{z} = [sW, vW, hW, dW]^\top$  in rapid task switching.

1296 The circuit receives four different inputs throughout each trial, which has a total length of 1.8s.

$$\mathbf{h} = \mathbf{h}_{\text{constant}} + \mathbf{h}_{\text{P,bias}} + \mathbf{h}_{\text{rule}} + \mathbf{h}_{\text{choice-period}} + \mathbf{h}_{\text{light}}. \quad (86)$$

1297 There is a constant input to every population,

$$\mathbf{h}_{\text{constant}} = I_{\text{constant}}[1, 1, 1, 1]^\top, \quad (87)$$

1298 a bias to the Pro populations

$$\mathbf{h}_{\text{P,bias}} = I_{\text{P,bias}}[1, 0, 1, 0]^\top, \quad (88)$$

1299 rule-based input depending on the condition

$$\mathbf{h}_{\text{P,rule}}(t) = \begin{cases} I_{\text{P,rule}}[1, 0, 1, 0]^\top, & \text{if } t \leq 1.2s \\ 0, & \text{otherwise} \end{cases} \quad (89)$$

1300

$$\mathbf{h}_{\text{A,rule}}(t) = \begin{cases} I_{\text{A,rule}}[0, 1, 0, 1]^\top, & \text{if } t \leq 1.2s \\ 0, & \text{otherwise} \end{cases}, \quad (90)$$

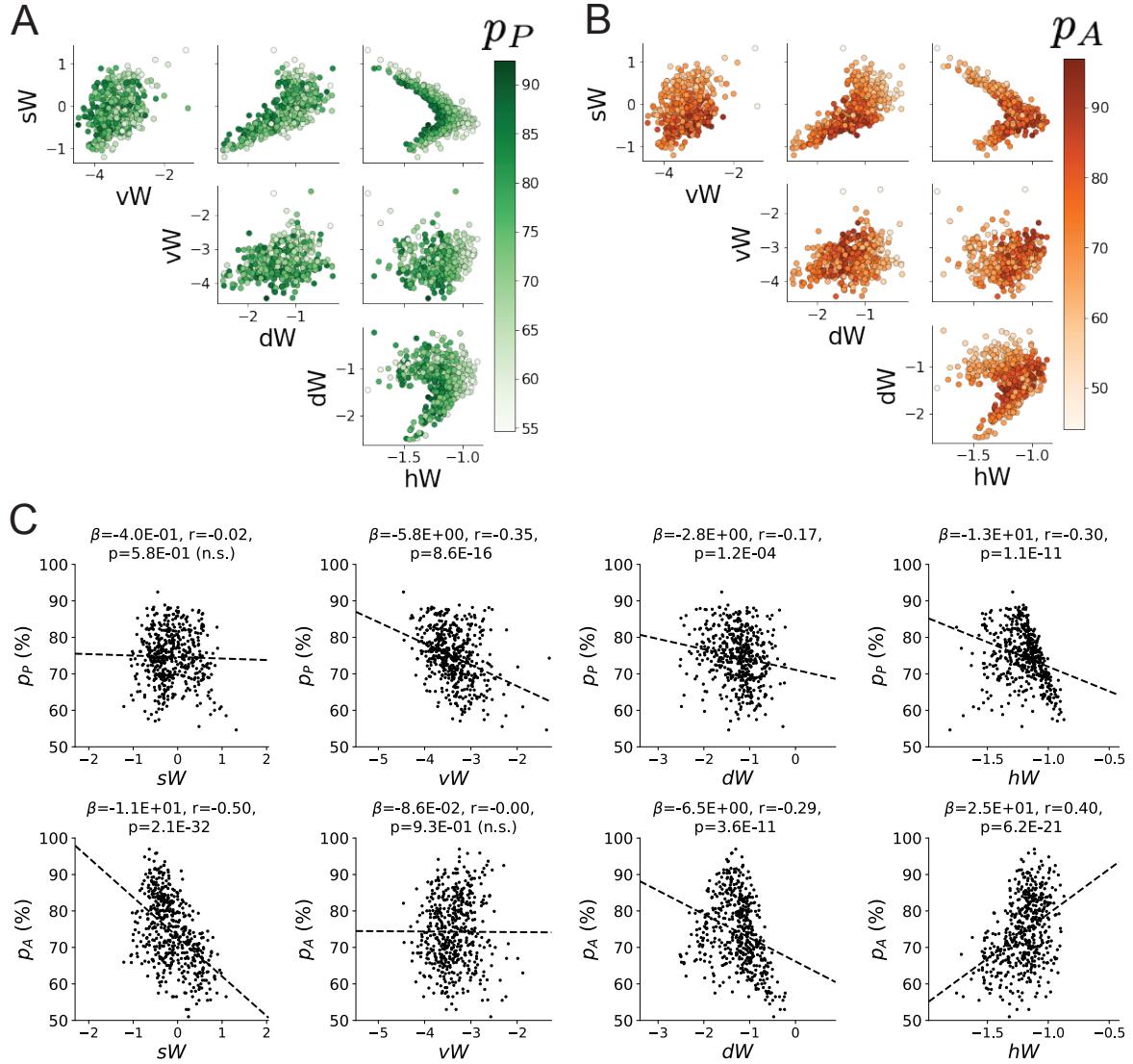


Figure 17: (SC1): **A.** Same pairplot as Fig. 4C colored by Pro task accuracy. **B.** Same as A colored by Anti task accuracy. **C.** Connectivity parameters of EPI distributions versus task accuracies.  $\beta$  is slope coefficient of linear regression,  $r$  is correlation, and  $p$  is the two-tailed p-value.

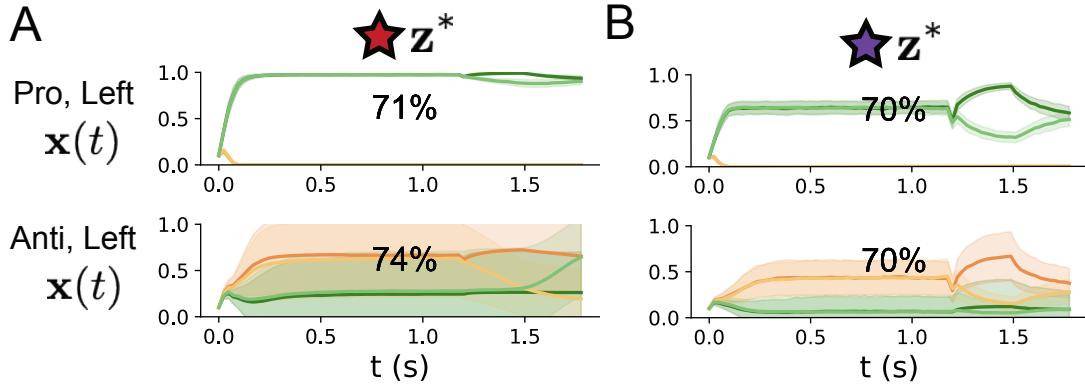


Figure 18: (SC2): **A.** Simulations in network regime 1 ( $hW_{\text{fixed}} = -1.5$ ). **B.** Simulations in network regime 2 ( $hW_{\text{fixed}} = -1.5$ ) .

1301 a choice-period input

$$\mathbf{h}_{\text{choice}}(t) = \begin{cases} I_{\text{choice}}[1, 1, 1, 1]^{\top}, & \text{if } t > 1.2s \\ 0, & \text{otherwise} \end{cases}, \quad (91)$$

1302 and an input to the right or left-side depending on where the light stimulus is delivered

$$\mathbf{h}_{\text{light}}(t) = \begin{cases} I_{\text{light}}[1, 1, 0, 0]^{\top}, & \text{if } 1.2s < t < 1.5s \text{ and Left} \\ I_{\text{light}}[0, 0, 1, 1]^{\top}, & \text{if } 1.2s < t < 1.5s \text{ and Right} \\ 0, & \text{otherwise} \end{cases}. \quad (92)$$

1303 The input parameterization was fixed to  $I_{\text{constant}} = 0.75$ ,  $I_{\text{P,bias}} = 0.5$ ,  $I_{\text{P,rule}} = 0.6$ ,  $I_{\text{A,rule}} = 0.6$ ,

1304  $I_{\text{choice}} = 0.25$ , and  $I_{\text{light}} = 0.5$ .

### 1305 5.5.2 Task accuracy calculation

1306 The accuracies of each task  $p_P$  and  $p_A$  are calculated as

$$p_P(\mathbf{x}; \mathbf{z}) = \mathbb{E}_{\mathbf{x}} [\Theta[x_{LP}(t = 1.8s) - x_{RP}(t = 1.8s)]] \quad (93)$$

1307 and

$$p_A(\mathbf{x}; \mathbf{z}) = \mathbb{E}_{\mathbf{x}} [\Theta[x_{RP}(t = 1.8s) - x_{LP}(t = 1.8s)]] \quad (94)$$

1308 given that the stimulus is on the left side, where  $\Theta$  is the Heaviside step function, and the accuracy  
1309 is averaged over 200 independent trials. The Heaviside step function is approximated as

$$\Theta(\mathbf{x}) = \text{sigmoid}(\beta \mathbf{x}), \quad (95)$$

1310 where  $\beta = 100$ .

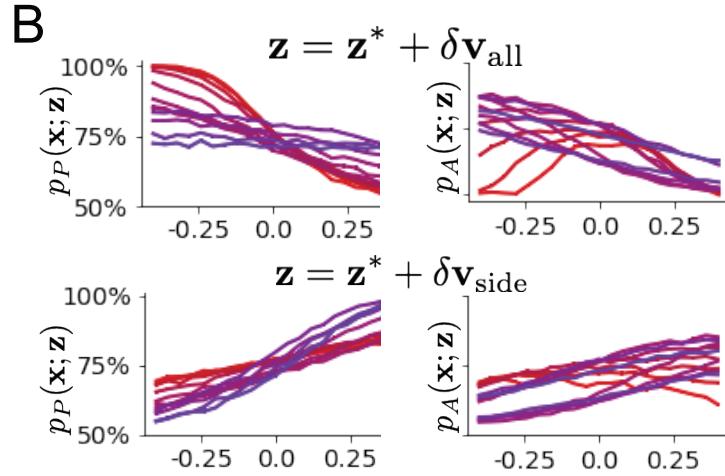
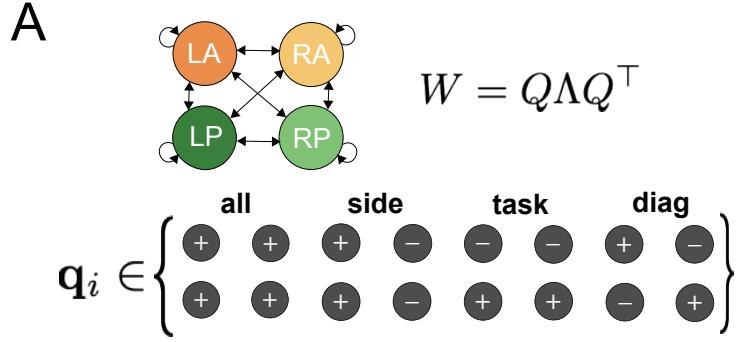


Figure 19: (SC3): **A.** Invariant eigenvectors of connectivity matrix  $W$ . **B.** Accuracies for connectivity perturbations for increasing  $\lambda_{\text{all}}$  and  $\lambda_{\text{side}}$  (rest shown in Fig. 4D).

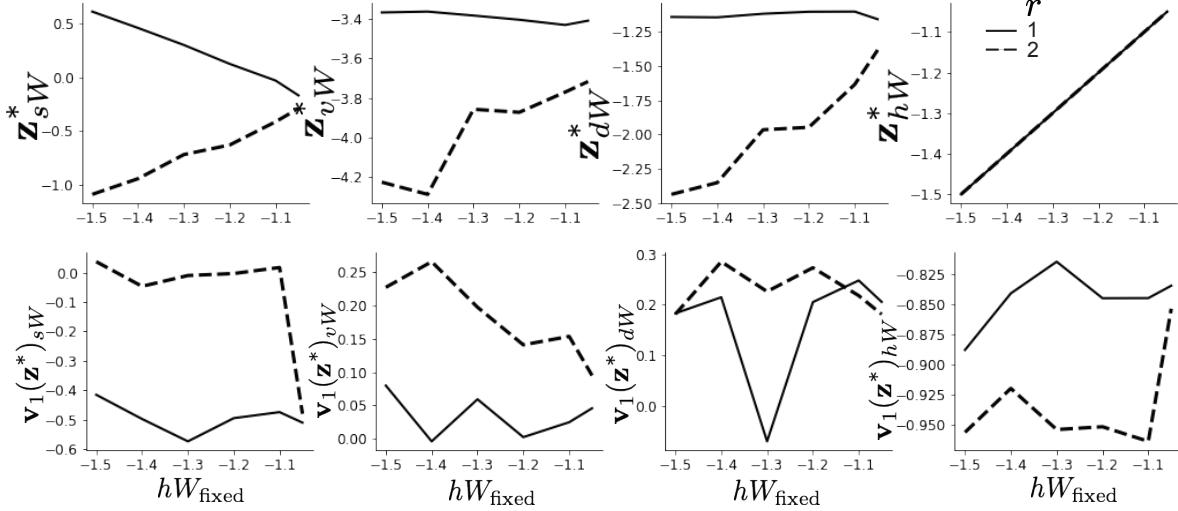


Figure 20: (SC4): **A.** The individual parameters of each mode throughout the two regimes. **B.** The individual sensitivities of parameters of each mode throughout the two regimes.

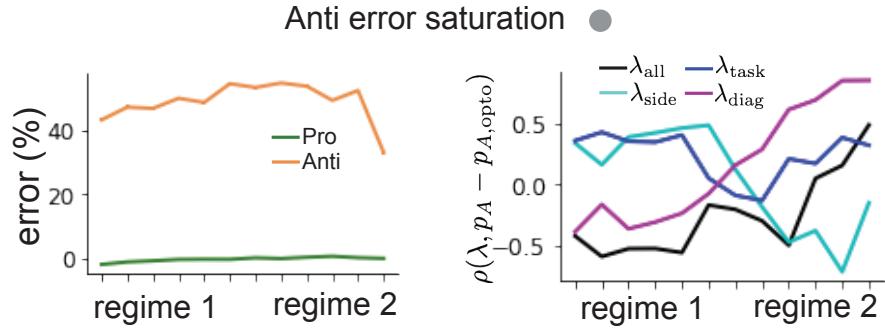


Figure 21: (SC5): (Left) Mean and standard error of Pro and Anti error from regime 1 to regime 2 at  $\gamma = 0.85$ . (Right) Correlations of connectivity eigenvalues with Anti error from regime 1 to regime 2 at  $\gamma = 0.85$ .

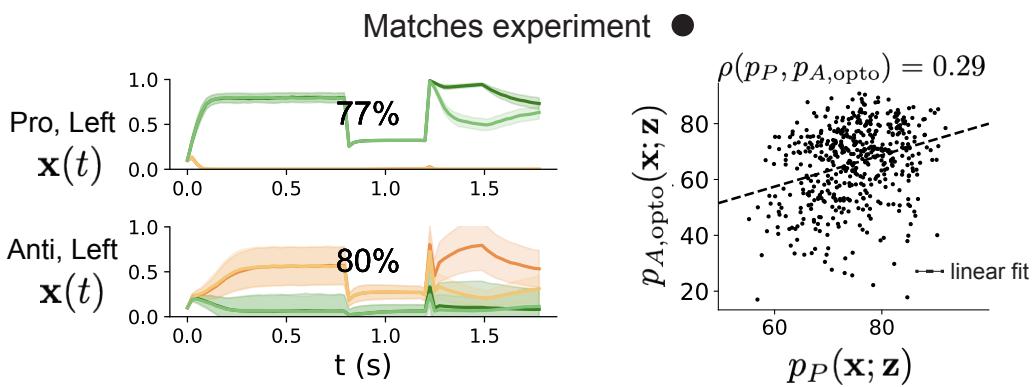


Figure 22: (SC6): (Left) Mean and standard deviation (shading) of responses of the SC model at the mode of the EPI distribution to delay period inactivation at  $\gamma = 0.675$ . (Right) Anti accuracy following delay period inactivation at  $\gamma = 0.675$  versus accuracy in the Pro task across connectivities in the EPI distribution.

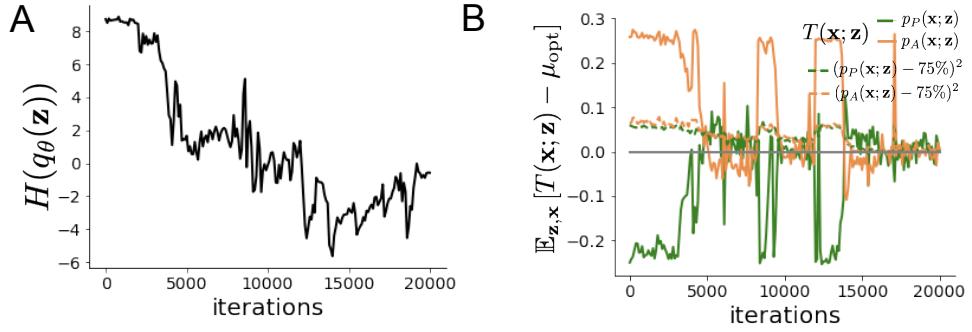


Figure 23: (SC7): **A.** Entropy throughout optimization. **B.** The emergent property statistic means and variances converge to their constraints at 20,000 iterations following the tenth augmented Lagrangian epoch.

1311 **5.5.3 EPI details for the SC model**

1312 Writing the EPI distribution as a maximum entropy distribution,  $T(\mathbf{x}; \mathbf{z})$  is comprised of both these  
1313 first and second moments of the accuracy in each task (as in Equations 17 and 18)

$$T(\mathbf{x}; \mathbf{z}) = \begin{bmatrix} p_P(\mathbf{x}; \mathbf{z}) \\ p_A(\mathbf{x}; \mathbf{z}) \\ (p_P(\mathbf{x}; \mathbf{z}) - .75)^2 \\ (p_A(\mathbf{x}; \mathbf{z}) - .75)^2 \end{bmatrix}, \quad (96)$$

$$\boldsymbol{\mu}_{\text{opt}} = \begin{bmatrix} .75 \\ .75 \\ .075^2 \\ .075^2 \end{bmatrix}. \quad (97)$$

1314 Throughout optimization, the augmented Lagrangian parameters  $\eta$  and  $c$ , were updated after each  
1315 epoch of 2,000 iterations (see Section 5.1.4). The optimization converged after ten epochs (Fig.  
1316 22).

1317 For EPI in Fig. 4C, we used a real NVP architecture with three coupling layers of affine transfor-  
1318 mations parameterized by two-layer neural networks of 50 units per layer. The initial distribution  
1319 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  
1320 augmented Lagrangian coefficient of  $c_0 = 10^2$ , a batch size  $n = 100$ , and  $\beta = 2$ . The distribution  
1321 was the greatest EPI distribution to converge across 5 random seeds with criteria  $N_{\text{test}} = 25$ .

1322 The bend in the EPI distribution is not a spurious result of the EPI optimization. The structure  
1323 discovered by EPI matches the shape of the set of points returned from brute-force random sampling

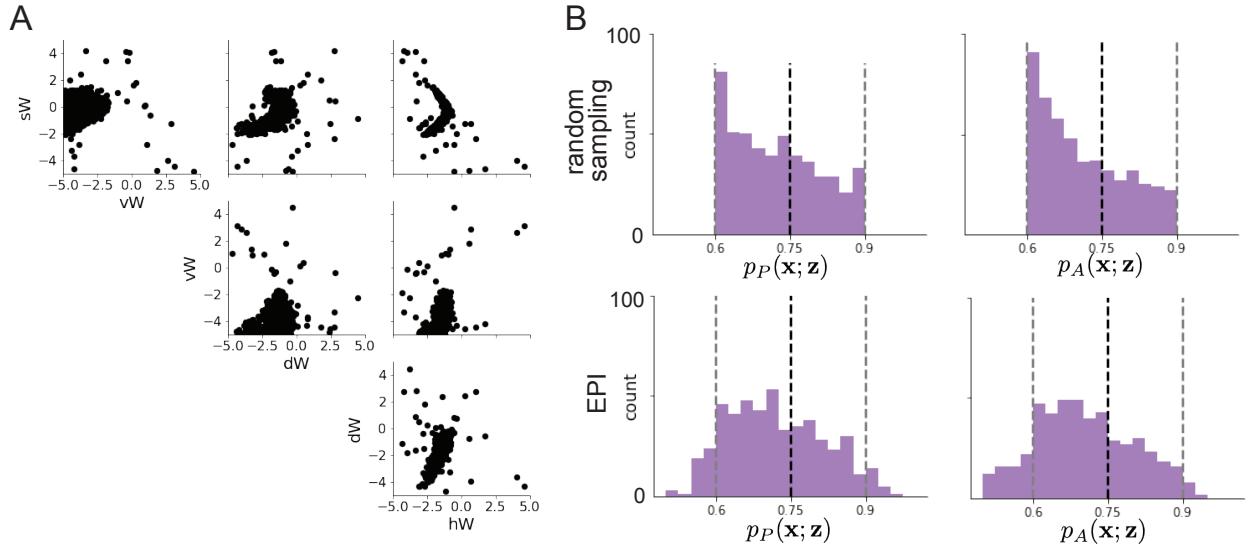


Figure 24: (SC8): **A.** Entropy throughout optimization. **B.** The emergent property statistic means and variances converge to their constraints at 20,000 iterations following the tenth augmented Lagrangian epoch.

(Fig. 24A) These connectivities were sampled from a uniform distribution over the range of each connectivity parameter, and all parameters producing accuracy in each task within the range of 60% to 90% were kept. This set of connectivities will not match the distribution of EPI exactly, since it is not conditioned on the emergent property. For example the parameter set returned by the brute-force search is biased towards lower accuracies (Fig. 24B).

#### 5.5.4 Regime identification with EPI

We sought two sets of parameters from  $q_{\theta}(\mathbf{z} | \mathcal{X})$  that were representative of each regime, so that we could assess their implications on computation. For fixed values of  $hW$ , we hypothesized that there are two modes: one in each regime of greater and lesser  $sw$ . To begin, we found one mode for each regime at  $hW_{\text{fixed}} = -1.5$  using 200 steps of gradient ascent of the deep probability distribution  $q_{\theta}(\mathbf{z} | \mathcal{X})$ . In regime 1, the initialization had positive  $sw$ , and the initialization had negative  $sw$  in regime 2, which led to disparate modes (Fig. 20 top). These modes were then used as the initialization to find the next mode at  $hW_{\text{fixed}} = -1.4$  and so on. 200 steps of gradient ascent were always taken, and learning rates of  $2.5 \times 10^{-4}$  and  $5 \times 10^{-4}$  were used for regimes 1 and 2, respectively. Each of these modes is denoted  $\mathbf{z}^*(hW_{\text{fixed}}, r)$  for regime  $r \in \{1, 2\}$ .

For the analyses in Figure 5C and Figure 21, we obtained parameters for each step along the continuum between regimes 1 and 2 by sampling from the EPI distribution. Each sample was

1342 assigned to the closest mode  $\mathbf{z}^*(hW_{\text{fixed}}, r)$ . Sampling continued until 500 samples were assigned to  
1343 each mode, which took 7.36 seconds. To obtain this many samples for each mode with brute force  
1344 sampling over the chosen prior, this would take 4.20 days.

1345 **5.5.5 Sensitivity analysis**

1346 At each mode, we measure the sensitivity dimension (that of most negative eigenvalue in the Hessian  
1347 of the EPI distribution)  $\mathbf{v}_1(\mathbf{z}^*)$ . To resolve sign degeneracy in eigenvectors, we chose  $\mathbf{v}_1(\mathbf{z}^*)$  to have  
1348 negative element in  $hW$ . This tells us what parameter combination rapid task switching is most  
1349 sensitive to at this parameter choice in the regime. We see that while the modes of each regime  
1350 gradually converge to similar connectivities at  $hW_{\text{fixed}} = -1.05$  (Fig. 20 top), the sensitivity  
1351 dimensions remain categorically different throughout the two regimes (Fig. 20 bottom). Only at  
1352  $hW_{\text{fixed}} = -1.05$  is there a flip in sensitivity from regime 2 to regime 1 (in  $\mathbf{v}_1(\mathbf{z}^*)_{sW}$  and  $\mathbf{v}_1(\mathbf{z}^*)_{hW}$ ).  
1353 There is thus some ambiguity regarding the “regime” of  $\mathbf{z}^*(-1.05, 2)$ , since the mode is derived  
1354 from an initialization in regime 2, but has sensitivity like regime 1. We can consider this as an  
1355 intermediate transitional region of parameter space between the two regimes. To emphasize this,  
1356  $\mathbf{z}^*(-1.05, 1)$  and  $\mathbf{z}^*(-1.05, 2)$  have the same color.

1357 **5.5.6 Connectivity eigendecomposition and processing modes**

1358 To understand the connectivity mechanisms governing task accuracy, we took the eigendecomposi-  
1359 tion of the symmetric connectivity matrices  $W = Q\Lambda Q^{-1}$ , which results in the same basis vectors  
1360  $\mathbf{q}_i$  for all  $W$  parameterized by  $\mathbf{z}$  (Fig. 19A). These basis vectors have intuitive roles in processing for  
1361 this task, and are accordingly named the *all* eigenmode - all neurons co-fluctuate, *side* eigenmode  
1362 - one side dominates the other, *task* eigenmode - the Pro or Anti populations dominate the other,  
1363 and *diag* mode - Pro- and Anti-populations of opposite hemispheres dominate the opposite pair.  
1364 Due to the parametric structure of the connectivity matrix, the parameters  $\mathbf{z}$  are a linear function  
1365 of the eigenvalues  $\boldsymbol{\lambda} = [\lambda_{\text{all}}, \lambda_{\text{side}}, \lambda_{\text{task}}, \lambda_{\text{diag}}]^T$  associated with these eigenmodes.

$$\mathbf{z} = A\boldsymbol{\lambda} \tag{98}$$

1366

$$A = \frac{1}{4} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & -1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & 1 & -1 \end{bmatrix}. \tag{99}$$

1367 We are interested in the effect of raising or lowering the amplification of each eigenmode in the  
1368 connectivity matrix. To test this, we calculate the unit vector of changes in the connectivity  $\mathbf{z}$  that  
1369 result from a change in the associated eigenvalues

$$\mathbf{v}_a = \frac{\frac{\partial \mathbf{z}}{\partial \lambda_a}}{\left\| \frac{\partial \mathbf{z}}{\partial \lambda_a} \right\|_2}, \quad (100)$$

1370 where

$$\frac{\partial \mathbf{z}}{\partial \lambda_a} = A \mathbf{e}_a, \quad (101)$$

1371 and e.g.  $\mathbf{e}_{\text{all}} = [1, 0, 0, 0]^\top$ . So  $\mathbf{v}_a$  is the normalized column of  $A$  corresponding to eigenmode  $a$ .  
1372 While perturbations in the sensitivity dimension  $\mathbf{v}_1(\mathbf{z}^*)$  adapt with the mode  $\mathbf{z}^*$  chosen, perturba-  
1373 tions in  $\mathbf{v}_a$  for  $a \in \{\text{all, side, text, diag}\}$  are invariant to  $\mathbf{z}$  (Equation 101).

1374 To understand the connectivity mechanisms that distinguish these two regimes, we perturb connec-  
1375 tivity at each mode in dimensions that have well defined roles in processing for the Pro and Anti  
1376 tasks. A convenient property of this connectivity parameterization is that there are  $\mathbf{z}$ -invariant  
1377 eigenmodes of connectivity, whose eigenvalues (or degree of amplification) change with  $\mathbf{z}$ . These  
1378 eigenmodes have intuitive roles in processing in each task, and are accordingly named the *all*,  
1379 *side*, *task*, and *diag* eigenmodes (see Section 5.5). Furthermore, the parameter dimension  $\mathbf{v}_a$   
1380 ( $a \in \{\text{all, side, task, and diag}\}$ ) that increases the eigenvalue of connectivity  $\lambda_a$  is  $\mathbf{z}$ -invariant (un-  
1381 like the sensitivity dimension  $\mathbf{v}_1(\mathbf{z})$ ) and  $\mathbf{v}_a \perp \mathbf{v}_{b \neq a}$ . Thus, by changing the degree of amplification  
1382 of each processing mode by perturbing  $\mathbf{z}$  along  $\mathbf{v}_a$ , we can elicit the differentiating properties of  
1383 the two regimes.

### 1384 5.5.7 Modeling optogenetic silencing.

1385 We tested whether the inferred SC model connectivities could reproduce experimental effects of  
1386 optogenetic inactivation in rats [69]. During periods of simulated optogenetic inactivation, activity  
1387 was decreased proportional to the optogenetic strength  $\gamma \in [0, 1]$

$$x_\alpha = (1 - \gamma)\phi(u_\alpha). \quad (102)$$

1388 Delay period inactivation was from  $0.8 < t < 1.2$ .