

Interrogating theoretical models of neural computation with deep inference
Sean R. Bittner¹, Agostina Palmigiano¹, Alex T. Piet^{2,3,4}, Chunyu A. Duan⁵, Carlos D. Brody^{2,3,6},
Kenneth D. Miller¹, and John P. Cunningham⁷.

¹Department of Neuroscience, Columbia University,

²Princeton Neuroscience Institute,

³Princeton University,

⁴Allen Institute for Brain Science,

⁵Institute of Neuroscience, Chinese Academy of Sciences,

⁶Howard Hughes Medical Institute,

⁷Department of Statistics, Columbia University

¹ 1 Abstract

² A cornerstone of theoretical neuroscience is the circuit model: a system of equations that captures a
³ hypothesized neural mechanism. Such models are valuable when they give rise to an experimentally
⁴ observed phenomenon – whether behavioral or in terms of neural activity – and thus can offer
⁵ insights into neural computation. The operation of these circuits, like all models, critically depends
⁶ on the choices of model parameters. When analytic derivation of the relationship between model
⁷ parameters and computational properties is intractable, approximate inference and simulation-
⁸ based techniques are relied upon for scientific insight. We bring the use of deep generative models
⁹ for probabilistic inference to bear on this problem, learning complex distributions of parameters
¹⁰ that produce the specified properties of computation. Our novel method solves the inverse problem
¹¹ by identifying the full space of parameters producing the emergent property. We motivate this
¹² methodology with a worked example analyzing sensitivity in the stomatogastric ganglion. We then
¹³ use it to reveal the key factors of variability in a model of primary visual cortex, gain a mechanistic
¹⁴ understanding of rapid task switching in superior colliculus models, and scale inference of large
¹⁵ low-rank RNN’s exhibiting stable amplification. This work illustrates how we can further leverage
¹⁶ the power of deep learning towards solving inverse problems in theoretical neuroscience.

₁₇ **2 Introduction**

₁₈ The fundamental practice of theoretical neuroscience is to use a mathematical model to understand
₁₉ neural computation, whether that computation enables perception, action, or some intermediate
₂₀ processing. A neural computation is systematized with a set of equations – the model – and
₂₁ these equations are motivated by biophysics, neurophysiology, and other conceptual considerations
₂₂ [1, 2, 3, 4]. The function of this system is governed by the choice of model *parameters*, which when
₂₃ configured in a particular way, give rise to a measurable signature of a computation. The work
₂₄ of analyzing a model then requires solving the inverse problem: given a computation of interest,
₂₅ how can we reason about particular parameter configurations? The inverse problem is crucial for
₂₆ reasoning about likely parameter values, uniquenesses and degeneracies, and predictions made by
₂₇ the model [5, 6].

₂₈ Consider the idealized practice: one carefully designs a model and analytically derives how com-
₂₉ putational properties determine model parameters. Seminal examples of this gold standard (which
₃₀ often adopt approaches from statistical physics) include our field’s understanding of memory ca-
₃₁ pacity in associative neural networks [7], chaos and autocorrelation timescales in random neural
₃₂ networks [8], the paradoxical effect [9], and decision making [10]. Unfortunately, as circuit models
₃₃ include more biological realism, theory via analytical derivation becomes intractable. Alternatively,
₃₄ we can gain insight into these complex models by identifying the full distribution of parameters con-
₃₅ sistent with specified emergent phenomena. By solving the inverse problem in this way, scientists
₃₆ can reason about the sensitivity and robustness of the model with respect to different parameter
₃₇ combinations [11, 12, 13, 6, 14].

₃₈ The preferred formalism for parameter identification in science is statistical inference, which has
₃₉ been used to great success in neuroscience through the stipulation of statistical generative models
₄₀ [15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29] (see review, [30]). However, most neural
₄₁ circuit models in theoretical neuroscience stipulate a noisy system of differential equations that can
₄₂ only be sampled or realized through forward simulation; they lack the explicit likelihood central to
₄₃ the probabilistic modeling toolkit. Therefore, the most popular approaches to the inverse problem
₄₄ have been likelihood-free methods such as approximate Bayesian computation (ABC) [31, 32], in
₄₅ which reasonable parameters are obtained via simulation and rejection.

₄₆ Of course, the challenge of doing inference in complex models has arisen in many scientific fields.
₄₇ In response, the machine learning community has made remarkable progress in recent years, via

48 the use of deep neural networks as powerful inference engines: a flexible function family that can
49 map observations back to probability distributions quantifying the likely parameter configurations.
50 One celebrated example of this approach from machine learning, of which we draw key inspiration
51 for this work, is the variational autoencoder (VAE) [33, 34], which uses a deep neural network
52 to induce an (approximate) posterior distribution on hidden variables in a latent variable model,
53 given data. Indeed, these tools have been used to great success in neuroscience as well, in particular
54 for interrogating hidden states in models of both cortical population activity [35, 36, 37, 38] and
55 animal behavior [39, 40, 41]. These works have used deep neural networks to expand the domain
56 of neural data sets amenable to statistical modeling [30].

57 Existing approaches to the inverse problem in theoretical neuroscience fall short in three key ways.
58 First, theoretical models of neural computation aim to reflect a complex biological reality, and as
59 a result, such models lack tractable likelihoods. Without an efficient calculation of the probability
60 of model properties given model parameters, neuroscientists resort to approximate Bayesian com-
61 putation [42, 43, 31], which requires a rejection heuristic, scales poorly, and only produces sets of
62 accepted parameters lacking probabilities. Second, there is an undesirable trade-off between the
63 flexibility and sampling speed of approximated posterior distributions. Sampling-based inference
64 approaches (e.g. ABC and Markov chain Monte Carlo (MCMC) [44, 45]) confer flexible approxima-
65 tions, yet scale poorly in number of parameters. While variational inference (VI) [46] often results
66 in fast posterior sampling, existing practice relies heavily on simplified classes of distributions [47].
67 Third, such parameter inference methods are designed to operate on experimentally collected data-
68 sets. Ultimately, the objects of interest in theoretical neuroscience are phenomena or features of
69 the model rather than singular data-sets.

70 To address these three challenges, we developed an inference methodology – ‘emergent property
71 inference’ – which learns a distribution over parameter configurations in a theoretical model. This
72 distribution has two critical properties: *(i)* it is chosen such that draws from the distribution (pa-
73 rameter configurations) correspond to systems of equations that give rise to a specified emergent
74 property (a set of constraints); and *(ii)* it is chosen to have maximum entropy given those con-
75 straints, such that we identify all likely parameters and can use the distribution to reason about
76 parametric sensitivity and degeneracies [48]. First, we use stochastic gradient techniques in the
77 spirit of likelihood-free variational inference [49] to enable inference in likelihood-free models of neu-
78 ral computation. Second, we stipulate a bijective deep neural network that induces a flexible family
79 of probability distributions over model parameterizations with a probability density we can calcu-

80 late [47, 50, 51], which confers fast sampling and sensitivity measurements. Third, we quantify the
81 notion of emergent properties as a set of moment constraints on datasets generated by the model.
82 Thus, an emergent property is not a single data realization, but a phenomenon or a feature of the
83 model. Conditioning on an emergent property requires a variant of deep probabilistic inference
84 methods, which we have previously introduced [52]. Taken together, emergent property inference
85 (EPI) provides a methodology for inferring parameter configurations consistent with a particular
86 emergent phenomena in theoretical models. We use a classic example of parametric degeneracy in
87 a biological system, the stomatogastric ganglion [53], to motivate and clarify the technical details
88 of EPI.

89 Equipped with this methodology, we then investigated three models of current importance in the-
90 oretical neuroscience. These models were chosen to demonstrate generality through ranges of bi-
91 ological realism (from conductance-based biophysics to recurrent neural networks), neural system
92 function (from pattern generation to decision making), and network scale (from four to hundreds of
93 neurons). First, we use EPI to understand the characteristics of noise across multiple neuron-type
94 populations that govern variability in a model of primary visual cortex. Then, we use EPI to infer
95 multiple regimes of superior colliculus connectivity that perform rapid task switching. The novel
96 scientific insights offered by EPI contextualize and clarify the previous studies exploring these mod-
97 els [54, 55]. Finally, we emphasize the scalability of EPI by inferring high-dimensional distributions
98 of RNNs exhibiting stable amplification. These results point to the value of deep inference for the
99 interrogation of biologically relevant models.

100 3 Results

101 3.1 Motivating emergent property inference of theoretical models

102 Consideration of the typical workflow of theoretical modeling clarifies the need for emergent prop-
103 erty inference. First, one designs or chooses an existing model that, it is hypothesized, captures the
104 computation of interest. To ground this process in a well-known example, consider the stomatoga-
105 stric ganglion (STG) of crustaceans, a small neural circuit which generates multiple rhythmic muscle
106 activation patterns for digestion [56]. Despite full knowledge of STG connectivity and a precise
107 characterization of its rhythmic pattern generation, biophysical models of the STG have compli-
108 cated relationships between circuit parameters and neural activity [53, 12]. A subcircuit model of
109 the STG [57] is shown schematically in Figure 1Emergent property inference (EPI) in the stomatoga-

tric ganglion. **A.** Conductance-based biophysical model of the STG subcircuit. In the STG model, jagged connections indicate electrical coupling having electrical conductance g_{el} . Other connections in the diagram are inhibitory synaptic projections having strength g_{synA} onto the hub neuron, and $g_{synB} = 5\text{nS}$ for mutual inhibitory connections. Parameters are represented by the vector \mathbf{z} and membrane potentials by the vector \mathbf{x} . The evolution of this model's activity $\mathbf{x}(t)$ is predicated by differential equations. **B.** Spiking frequency $\omega(\mathbf{x}; \mathbf{z})$ is an emergent property statistic. In this example, spiking frequency is measured from simulated activity of the STG model at parameter choices of $g_{el} = 4.5\text{nS}$ and $g_{synA} = 3\text{nS}$. **C.** The emergent property of intermediate hub frequency, in which the hub neuron fires at a rate between the fast and slow frequencies. This emergent property is defined by a mean and variance on the emergent property statistic. Simulated activity traces are colored by log probability density of their generating parameters in the EPI-inferred distribution (Panel E). **D.** For a choice of model and emergent property, emergent property inference (EPI) learns a deep probability distribution of parameters \mathbf{z} . Deep probability distributions map a simple random variable \mathbf{z}_0 through a deep neural network with weights and biases $\boldsymbol{\theta}$ to parameters $\mathbf{z} = g_{\boldsymbol{\theta}}(\mathbf{z}_0)$. In EPI optimization, stochastic gradient steps in $\boldsymbol{\theta}$ are taken such that entropy is maximized, and the emergent property \mathcal{X} is produced. The EPI posterior distribution is denoted $q_{\boldsymbol{\theta}}(\mathbf{z} | \mathcal{X})$. **E.** The EPI posterior producing intermediate hub frequency. Samples are colored by log probability density. Distribution contours of average hub neuron frequency from mean of .55 Hz are shown at levels of .525, .53,575 Hz (dark to light gray away from mean). Eigenvectors of the Hessian at the mode of the inferred distribution are indicated as \mathbf{v}_1 (solid) and \mathbf{v}_2 (dashed) with lengths scaled by the square root of the absolute value of their eigenvalues. **F** Simulations from parameters in E. (Top) The predictive distribution of the posterior obeys the emergent property. The black and gray dashed lines show the mean and two standard deviations according the emergent property, respectively. (Bottom) Simulations at the starred parameter valuesfigure.1A, and note that the behavior of this model will be critically dependent on its parameterization – the choices of conductance parameters $\mathbf{z} = [g_{el}, g_{synA}]$. Specifically, the two fast neurons (f_1 and f_2) mutually inhibit one another, and oscillate at a faster frequency than the mutually inhibiting slow neurons (s_1 and s_2). The hub neuron (hub) couples with either the fast or slow population or both.

Second, once the model is selected, one defines the emergent phenomena of scientific interest. In the STG example, we are concerned with neural spiking frequency, which emerges from the dynamics of the circuit model 1Emergent property inference (EPI) in the stomatogastric ganglion. **A.** Conductance-based biophysical model of the STG subcircuit. In the STG model, jagged connections indicate electrical coupling having electrical conductance g_{el} . Other connections in the diagram are inhibitory synaptic projections having strength g_{synA} onto the hub neuron, and $g_{synB} = 5\text{nS}$ for mutual inhibitory connections.

142 Parameters are represented by the vector \mathbf{z} and membrane potentials by the vector \mathbf{x} . The evolution of this
 143 model's activity $\mathbf{x}(t)$ is predicated by differential equations. **B.** Spiking frequency $\omega(\mathbf{x}; \mathbf{z})$ is an emergent
 144 property statistic. In this example, spiking frequency is measured from simulated activity of the STG
 145 model at parameter choices of $g_{el} = 4.5\text{nS}$ and $g_{synA} = 3\text{nS}$. **C.** The emergent property of intermediate hub
 146 frequency, in which the hub neuron fires at a rate between the fast and slow frequencies. This emergent
 147 property is defined by a mean and variance on the emergent property statistic. Simulated activity traces are
 148 colored by log probability density of their generating parameters in the EPI-inferred distribution (Panel E).
 149 **D.** For a choice of model and emergent property, emergent property inference (EPI) learns a deep probability
 150 distribution of parameters \mathbf{z} . Deep probability distributions map a simple random variable \mathbf{z}_0 through a
 151 deep neural network with weights and biases $\boldsymbol{\theta}$ to parameters $\mathbf{z} = g_{\boldsymbol{\theta}}(\mathbf{z}_0)$. In EPI optimization, stochastic
 152 gradient steps in $\boldsymbol{\theta}$ are taken such that entropy is maximized, and the emergent property \mathcal{X} is produced.
 153 The EPI posterior distribution is denoted $q_{\boldsymbol{\theta}}(\mathbf{z} \mid \mathcal{X})$. **E.** The EPI posterior producing intermediate hub
 154 frequency. Samples are colored by log probability density. Distribution contours of average hub neuron
 155 frequency from mean of .55 Hz are shown at levels of .525, .53,575 Hz (dark to light gray away from
 156 mean). Eigenvectors of the Hessian at the mode of the inferred distribution are indicated as \mathbf{v}_1 (solid) and
 157 \mathbf{v}_2 (dashed) with lengths scaled by the square root of the absolute value of their eigenvalues. **F** Simulations
 158 from parameters in E. (Top) The predictive distribution of the posterior obeys the emergent property. The
 159 black and gray dashed lines show the mean and two standard deviations according the emergent property,
 160 respectively. (Bottom) Simulations at the starred parameter valuesfigure.1B. An interesting emergent
 161 property of this stochastic model is when the hub neuron fires at an intermediate frequency between
 162 the intrinsic spiking rates of the fast and slow populations. This emergent property is shown
 163 in Figure 1Emergent property inference (EPI) in the stomatogastric ganglion. **A.** Conductance-based
 164 biophysical model of the STG subcircuit. In the STG model, jagged connections indicate electrical coupling
 165 having electrical conductance g_{el} . Other connections in the diagram are inhibitory synaptic projections
 166 having strength g_{synA} onto the hub neuron, and $g_{synB} = 5\text{nS}$ for mutual inhibitory connections. Parameters
 167 are represented by the vector \mathbf{z} and membrane potentials by the vector \mathbf{x} . The evolution of this model's
 168 activity $\mathbf{x}(t)$ is predicated by differential equations. **B.** Spiking frequency $\omega(\mathbf{x}; \mathbf{z})$ is an emergent property
 169 statistic. In this example, spiking frequency is measured from simulated activity of the STG model at
 170 parameter choices of $g_{el} = 4.5\text{nS}$ and $g_{synA} = 3\text{nS}$. **C.** The emergent property of intermediate hub frequency,
 171 in which the hub neuron fires at a rate between the fast and slow frequencies. This emergent property is
 172 defined by a mean and variance on the emergent property statistic. Simulated activity traces are colored
 173 by log probability density of their generating parameters in the EPI-inferred distribution (Panel E). **D.**
 174 For a choice of model and emergent property, emergent property inference (EPI) learns a deep probability

175 distribution of parameters \mathbf{z} . Deep probability distributions map a simple random variable \mathbf{z}_0 through a
176 deep neural network with weights and biases $\boldsymbol{\theta}$ to parameters $\mathbf{z} = g_{\boldsymbol{\theta}}(\mathbf{z}_0)$. In EPI optimization, stochastic
177 gradient steps in $\boldsymbol{\theta}$ are taken such that entropy is maximized, and the emergent property \mathcal{X} is produced.
178 The EPI posterior distribution is denoted $q_{\boldsymbol{\theta}}(\mathbf{z} \mid \mathcal{X})$. **E.** The EPI posterior producing intermediate hub
179 frequency. Samples are colored by log probability density. Distribution contours of average hub neuron
180 frequency from mean of .55 Hz are shown at levels of .525, .53,575 Hz (dark to light gray away from
181 mean). Eigenvectors of the Hessian at the mode of the inferred distribution are indicated as \mathbf{v}_1 (solid) and
182 \mathbf{v}_2 (dashed) with lengths scaled by the square root of the absolute value of their eigenvalues. **F** Simulations
183 from parameters in E. (Top) The predictive distribution of the posterior obeys the emergent property. The
184 black and gray dashed lines show the mean and two standard deviations according the emergent property,
185 respectively. (Bottom) Simulations at the starred parameter valuesfigure.1C at an average frequency of
186 0.55Hz.

187 Third, parameter analyses ensue: brute-force parameter sweeps, ABC sampling, and sensitivity
188 analyses are all routinely used to reason about what parameter configurations lead to an emergent
189 property. In this last step lies the opportunity for a precise quantification of the emergent property
190 as a statistical feature of the model. Once we have such a methodology, we can infer a probability
191 distribution over parameter configurations that produce this emergent property.

192 Before presenting technical details (in the following section), let us understand emergent property
193 inference schematically: EPI (Fig. 1Emergent property inference (EPI) in the stomatogastric ganglion.
194 **A.** Conductance-based biophysical model of the STG subcircuit. In the STG model, jagged connections in-
195 dicate electrical coupling having electrical conductance g_{el} . Other connections in the diagram are inhibitory
196 synaptic projections having strength g_{synA} onto the hub neuron, and $g_{synB} = 5nS$ for mutual inhibitory
197 connections. Parameters are represented by the vector \mathbf{z} and membrane potentials by the vector \mathbf{x} . The
198 evolution of this model's activity $\mathbf{x}(t)$ is predicated by differential equations. **B.** Spiking frequency $\omega(\mathbf{x}; \mathbf{z})$
199 is an emergent property statistic. In this example, spiking frequency is measured from simulated activity of
200 the STG model at parameter choices of $g_{el} = 4.5nS$ and $g_{synA} = 3nS$. **C.** The emergent property of inter-
201 mediate hub frequency, in which the hub neuron fires at a rate between the fast and slow frequencies. This
202 emergent property is defined by a mean and variance on the emergent property statistic. Simulated activity
203 traces are colored by log probability density of their generating parameters in the EPI-inferred distribution
204 (Panel E). **D.** For a choice of model and emergent property, emergent property inference (EPI) learns a deep
205 probability distribution of parameters \mathbf{z} . Deep probability distributions map a simple random variable \mathbf{z}_0 through a deep neural network with weights and biases $\boldsymbol{\theta}$ to parameters $\mathbf{z} = g_{\boldsymbol{\theta}}(\mathbf{z}_0)$. In EPI optimization,
206

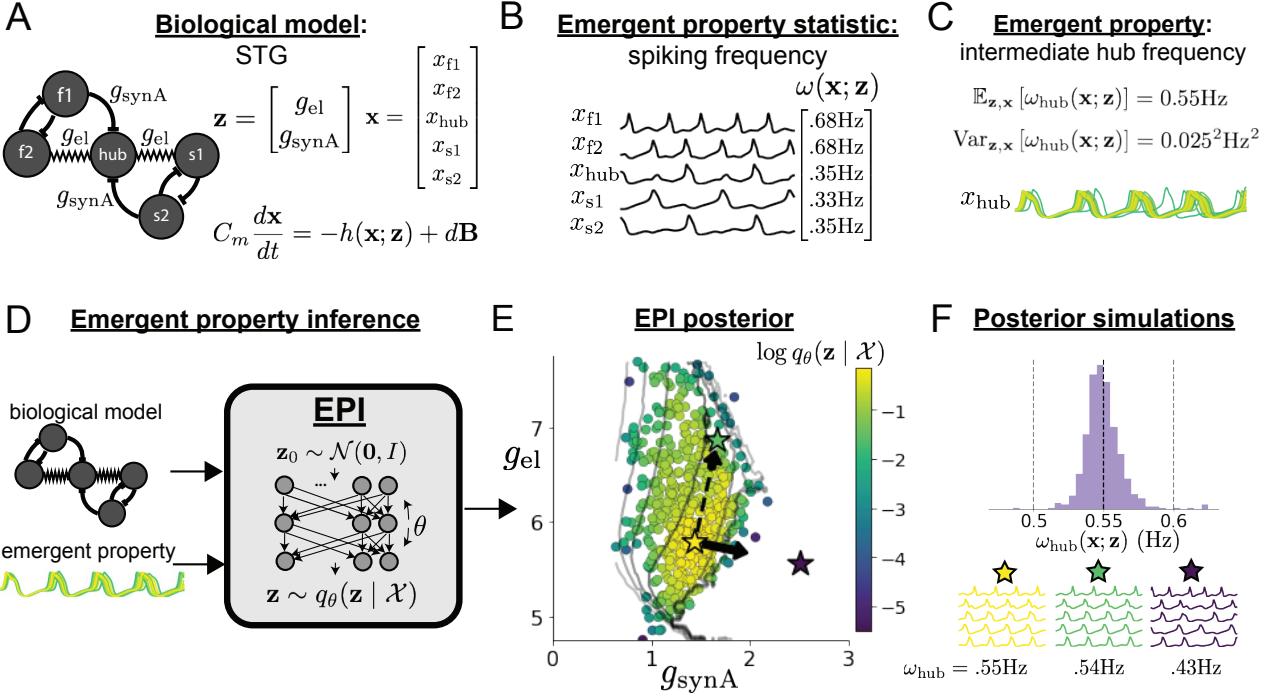


Figure 1: Emergent property inference (EPI) in the stomatogastric ganglion. **A.** Conductance-based biophysical model of the STG subcircuit. In the STG model, jagged connections indicate electrical coupling having electrical conductance g_{el} . Other connections in the diagram are inhibitory synaptic projections having strength g_{synA} onto the hub neuron, and $g_{synB} = 5\text{nS}$ for mutual inhibitory connections. Parameters are represented by the vector \mathbf{z} and membrane potentials by the vector \mathbf{x} . The evolution of this model's activity $\mathbf{x}(t)$ is predicated by differential equations. **B.** Spiking frequency $\omega(\mathbf{x}; \mathbf{z})$ is an emergent property statistic. In this example, spiking frequency is measured from simulated activity of the STG model at parameter choices of $g_{el} = 4.5\text{nS}$ and $g_{synA} = 3\text{nS}$. **C.** The emergent property of intermediate hub frequency, in which the hub neuron fires at a rate between the fast and slow frequencies. This emergent property is defined by a mean and variance on the emergent property statistic. Simulated activity traces are colored by log probability density of their generating parameters in the EPI-inferred distribution (Panel E). **D.** For a choice of model and emergent property, emergent property inference (EPI) learns a deep probability distribution of parameters \mathbf{z} . Deep probability distributions map a simple random variable $\mathbf{z}_0 \sim \mathcal{N}(0, I)$ through a deep neural network with weights and biases $\boldsymbol{\theta}$ to parameters $\mathbf{z} = q_{\boldsymbol{\theta}}(\mathbf{z}_0)$. In EPI optimization, stochastic gradient steps in $\boldsymbol{\theta}$ are taken such that entropy is maximized, and the emergent property \mathcal{X} is produced. The EPI posterior distribution is denoted $q_{\boldsymbol{\theta}}(\mathbf{z} | \mathcal{X})$. **E.** The EPI posterior producing intermediate hub frequency. Samples are colored by log probability density. Distribution contours of average hub neuron frequency from mean of .55 Hz are shown at levels of .525, .53,575 Hz (dark to light gray away from mean). Eigenvectors of the Hessian at the mode of the inferred distribution are indicated as \mathbf{v}_1 (solid) and \mathbf{v}_2 (dashed) with lengths scaled by the square root of the absolute value of their eigenvalues. **F** Simulations from parameters in E. (Top) The predictive distribution of the posterior obeys the emergent property. The black and gray dashed lines show the mean and two standard deviations according the emergent property, respectively. (Bottom) Simulations at the starred parameter values.

207 stochastic gradient steps in θ are taken such that entropy is maximized, and the emergent property \mathcal{X} is
 208 produced. The EPI posterior distribution is denoted $q_\theta(\mathbf{z} | \mathcal{X})$. **E.** The EPI posterior producing intermediate
 209 hub frequency. Samples are colored by log probability density. Distribution contours of average hub neuron
 210 frequency from mean of .55 Hz are shown at levels of .525, .53,575 Hz (dark to light gray away from
 211 mean). Eigenvectors of the Hessian at the mode of the inferred distribution are indicated as \mathbf{v}_1 (solid) and
 212 \mathbf{v}_2 (dashed) with lengths scaled by the square root of the absolute value of their eigenvalues. **F** Simulations
 213 from parameters in E. (Top) The predictive distribution of the posterior obeys the emergent property. The
 214 black and gray dashed lines show the mean and two standard deviations according the emergent property,
 215 respectively. (Bottom) Simulations at the starred parameter valuesfigure.1D) takes, as input, the model
 216 and the specified emergent property, and as its output, produces the parameter distribution EPI
 217 (Fig. 1E)Emergent property inference (EPI) in the stomatogastric ganglion. **A.** Conductance-based biophys-
 218 ical model of the STG subcircuit. In the STG model, jagged connections indicate electrical coupling having
 219 electrical conductance g_{el} . Other connections in the diagram are inhibitory synaptic projections having
 220 strength g_{synA} onto the hub neuron, and $g_{synB} = 5nS$ for mutual inhibitory connections. Parameters are
 221 represented by the vector \mathbf{z} and membrane potentials by the vector \mathbf{x} . The evolution of this model's activity
 222 $\mathbf{x}(t)$ is predicated by differential equations. **B.** Spiking frequency $\omega(\mathbf{x}; \mathbf{z})$ is an emergent property statistic.
 223 In this example, spiking frequency is measured from simulated activity of the STG model at parameter
 224 choices of $g_{el} = 4.5nS$ and $g_{synA} = 3nS$. **C.** The emergent property of intermediate hub frequency, in which
 225 the hub neuron fires at a rate between the fast and slow frequencies. This emergent property is defined
 226 by a mean and variance on the emergent property statistic. Simulated activity traces are colored by log
 227 probability density of their generating parameters in the EPI-inferred distribution (Panel E). **D.** For a choice
 228 of model and emergent property, emergent property inference (EPI) learns a deep probability distribution
 229 of parameters \mathbf{z} . Deep probability distributions map a simple random variable \mathbf{z}_0 through a deep neural
 230 network with weights and biases θ to parameters $\mathbf{z} = g_\theta(\mathbf{z}_0)$. In EPI optimization, stochastic gradient steps
 231 in θ are taken such that entropy is maximized, and the emergent property \mathcal{X} is produced. The EPI posterior
 232 distribution is denoted $q_\theta(\mathbf{z} | \mathcal{X})$. **E.** The EPI posterior producing intermediate hub frequency. Samples are
 233 colored by log probability density. Distribution contours of average hub neuron frequency from mean of .55
 234 Hz are shown at levels of .525, .53,575 Hz (dark to light gray away from mean). Eigenvectors of the
 235 Hessian at the mode of the inferred distribution are indicated as \mathbf{v}_1 (solid) and \mathbf{v}_2 (dashed) with lengths
 236 scaled by the square root of the absolute value of their eigenvalues. **F** Simulations from parameters in E.
 237 (Top) The predictive distribution of the posterior obeys the emergent property. The black and gray dashed
 238 lines show the mean and two standard deviations according the emergent property, respectively. (Bottom)
 239 Simulations at the starred parameter valuesfigure.1E). This distribution – represented for clarity as

240 samples from the distribution – is then a scientifically meaningful and mathematically tractable
 241 object. In the STG model, this distribution can be specifically queried to reveal the prototypical
 242 parameter configuration for network syncing (the mode; Figure 1E) emergent property inference (EPI) in
 243 the stomatogastric ganglion. **A.** Conductance-based biophysical model of the STG subcircuit. In the STG
 244 model, jagged connections indicate electrical coupling having electrical conductance g_{el} . Other connections in
 245 the diagram are inhibitory synaptic projections having strength g_{synA} onto the hub neuron, and $g_{synB} = 5\text{nS}$
 246 for mutual inhibitory connections. Parameters are represented by the vector \mathbf{z} and membrane potentials by
 247 the vector \mathbf{x} . The evolution of this model’s activity $\mathbf{x}(t)$ is predicated by differential equations. **B.** Spiking
 248 frequency $\omega(\mathbf{x}; \mathbf{z})$ is an emergent property statistic. In this example, spiking frequency is measured from
 249 simulated activity of the STG model at parameter choices of $g_{el} = 4.5\text{nS}$ and $g_{synA} = 3\text{nS}$. **C.** The emergent
 250 property of intermediate hub frequency, in which the hub neuron fires at a rate between the fast and slow
 251 frequencies. This emergent property is defined by a mean and variance on the emergent property statistic.
 252 Simulated activity traces are colored by log probability density of their generating parameters in the EPI-
 253 inferred distribution (Panel E). **D.** For a choice of model and emergent property, emergent property inference
 254 (EPI) learns a deep probability distribution of parameters \mathbf{z} . Deep probability distributions map a simple
 255 random variable \mathbf{z}_0 through a deep neural network with weights and biases $\boldsymbol{\theta}$ to parameters $\mathbf{z} = g_{\boldsymbol{\theta}}(\mathbf{z}_0)$. In
 256 EPI optimization, stochastic gradient steps in $\boldsymbol{\theta}$ are taken such that entropy is maximized, and the emergent
 257 property \mathcal{X} is produced. The EPI posterior distribution is denoted $q_{\boldsymbol{\theta}}(\mathbf{z} | \mathcal{X})$. **E.** The EPI posterior pro-
 258 ducing intermediate hub frequency. Samples are colored by log probability density. Distribution contours of
 259 average hub neuron frequency from mean of .55 Hz are shown at levels of .525, .53,575 Hz (dark to light
 260 gray away from mean). Eigenvectors of the Hessian at the mode of the inferred distribution are indicated
 261 as \mathbf{v}_1 (solid) and \mathbf{v}_2 (dashed) with lengths scaled by the square root of the absolute value of their eigen-
 262 values. **F** Simulations from parameters in E. (Top) The predictive distribution of the posterior obeys the
 263 emergent property. The black and gray dashed lines show the mean and two standard deviations according
 264 the emergent property, respectively. (Bottom) Simulations at the starred parameter values (figure 1E yellow
 265 star), and how network syncing decays based on changes away from the mode. The eigenvectors
 266 (of the Hessian of the distribution at the mode) quantitatively formalize the robustness of inter-
 267 mediate hub frequency (Fig. 1E) emergent property inference (EPI) in the stomatogastric ganglion. **A.**
 268 Conductance-based biophysical model of the STG subcircuit. In the STG model, jagged connections indi-
 269 cate electrical coupling having electrical conductance g_{el} . Other connections in the diagram are inhibitory
 270 synaptic projections having strength g_{synA} onto the hub neuron, and $g_{synB} = 5\text{nS}$ for mutual inhibitory
 271 connections. Parameters are represented by the vector \mathbf{z} and membrane potentials by the vector \mathbf{x} . The
 272 evolution of this model’s activity $\mathbf{x}(t)$ is predicated by differential equations. **B.** Spiking frequency $\omega(\mathbf{x}; \mathbf{z})$

273 is an emergent property statistic. In this example, spiking frequency is measured from simulated activity of
 274 the STG model at parameter choices of $g_{el} = 4.5\text{nS}$ and $g_{synA} = 3\text{nS}$. **C.** The emergent property of inter-
 275 mediate hub frequency, in which the hub neuron fires at a rate between the fast and slow frequencies. This
 276 emergent property is defined by a mean and variance on the emergent property statistic. Simulated activity
 277 traces are colored by log probability density of their generating parameters in the EPI-inferred distribution
 278 (Panel E). **D.** For a choice of model and emergent property, emergent property inference (EPI) learns a deep
 279 probability distribution of parameters \mathbf{z} . Deep probability distributions map a simple random variable \mathbf{z}_0
 280 through a deep neural network with weights and biases $\boldsymbol{\theta}$ to parameters $\mathbf{z} = g_{\boldsymbol{\theta}}(\mathbf{z}_0)$. In EPI optimization,
 281 stochastic gradient steps in $\boldsymbol{\theta}$ are taken such that entropy is maximized, and the emergent property \mathcal{X} is
 282 produced. The EPI posterior distribution is denoted $q_{\boldsymbol{\theta}}(\mathbf{z} \mid \mathcal{X})$. **E.** The EPI posterior producing intermediate
 283 hub frequency. Samples are colored by log probability density. Distribution contours of average hub neuron
 284 frequency from mean of .55 Hz are shown at levels of .525, .53,575 Hz (dark to light gray away from
 285 mean). Eigenvectors of the Hessian at the mode of the inferred distribution are indicated as \mathbf{v}_1 (solid) and
 286 \mathbf{v}_2 (dashed) with lengths scaled by the square root of the absolute value of their eigenvalues. **F** Simulations
 287 from parameters in E. (Top) The predictive distribution of the posterior obeys the emergent property. The
 288 black and gray dashed lines show the mean and two standard deviations according the emergent property,
 289 respectively. (Bottom) Simulations at the starred parameter valuesfigure.1E solid (v_1) and dashed (v_2)
 290 black arrows). Indeed, samples equidistant from the mode along these EPI-identified dimensions
 291 of sensitivity (v_1) and degeneracy (v_2) agree with error contours (Fig. 1Emergent property inference
 292 (EPI) in the stomatogastric ganglion. **A.** Conductance-based biophysical model of the STG subcircuit. In
 293 the STG model, jagged connections indicate electrical coupling having electrical conductance g_{el} . Other
 294 connections in the diagram are inhibitory synaptic projections having strength g_{synA} onto the hub neu-
 295 ron, and $g_{synB} = 5\text{nS}$ for mutual inhibitory connections. Parameters are represented by the vector \mathbf{z} and
 296 membrane potentials by the vector \mathbf{x} . The evolution of this model's activity $\mathbf{x}(t)$ is predicated by differ-
 297 ential equations. **B.** Spiking frequency $\omega(\mathbf{x}; \mathbf{z})$ is an emergent property statistic. In this example, spiking
 298 frequency is measured from simulated activity of the STG model at parameter choices of $g_{el} = 4.5\text{nS}$ and
 299 $g_{synA} = 3\text{nS}$. **C.** The emergent property of intermediate hub frequency, in which the hub neuron fires at
 300 a rate between the fast and slow frequencies. This emergent property is defined by a mean and variance
 301 on the emergent property statistic. Simulated activity traces are colored by log probability density of their
 302 generating parameters in the EPI-inferred distribution (Panel E). **D.** For a choice of model and emergent
 303 property, emergent property inference (EPI) learns a deep probability distribution of parameters \mathbf{z} . Deep
 304 probability distributions map a simple random variable \mathbf{z}_0 through a deep neural network with weights and
 305 biases $\boldsymbol{\theta}$ to parameters $\mathbf{z} = g_{\boldsymbol{\theta}}(\mathbf{z}_0)$. In EPI optimization, stochastic gradient steps in $\boldsymbol{\theta}$ are taken such that

306 entropy is maximized, and the emergent property \mathcal{X} is produced. The EPI posterior distribution is denoted
 307 $q_{\theta}(\mathbf{z} \mid \mathcal{X})$. **E.** The EPI posterior producing intermediate hub frequency. Samples are colored by log prob-
 308 ability density. Distribution contours of average hub neuron frequency from mean of .55 Hz are shown at
 309 levels of .525, .53,575 Hz (dark to light gray away from mean). Eigenvectors of the Hessian at the mode
 310 of the inferred distribution are indicated as \mathbf{v}_1 (solid) and \mathbf{v}_2 (dashed) with lengths scaled by the square
 311 root of the absolute value of their eigenvalues. **F** Simulations from parameters in E. (Top) The predictive
 312 distribution of the posterior obeys the emergent property. The black and gray dashed lines show the mean
 313 and two standard deviations according the emergent property, respectively. (Bottom) Simulations at the
 314 starred parameter values (figure 1E contours) and have diminished or preserved hub frequency, respec-
 315 tively (Fig. 1E). **A.** Emergent property inference (EPI) in the stomatogastric ganglion. **A.** Conductance-based
 316 biophysical model of the STG subcircuit. In the STG model, jagged connections indicate electrical coupling
 317 having electrical conductance g_{el} . Other connections in the diagram are inhibitory synaptic projections hav-
 318 ing strength g_{synA} onto the hub neuron, and $g_{synB} = 5nS$ for mutual inhibitory connections. Parameters are
 319 represented by the vector \mathbf{z} and membrane potentials by the vector \mathbf{x} . The evolution of this model's activity
 320 $\mathbf{x}(t)$ is predicated by differential equations. **B.** Spiking frequency $\omega(\mathbf{x}; \mathbf{z})$ is an emergent property statistic.
 321 In this example, spiking frequency is measured from simulated activity of the STG model at parameter
 322 choices of $g_{el} = 4.5nS$ and $g_{synA} = 3nS$. **C.** The emergent property of intermediate hub frequency, in which
 323 the hub neuron fires at a rate between the fast and slow frequencies. This emergent property is defined
 324 by a mean and variance on the emergent property statistic. Simulated activity traces are colored by log
 325 probability density of their generating parameters in the EPI-inferred distribution (Panel E). **D.** For a choice
 326 of model and emergent property, emergent property inference (EPI) learns a deep probability distribution
 327 of parameters \mathbf{z} . Deep probability distributions map a simple random variable \mathbf{z}_0 through a deep neural
 328 network with weights and biases θ to parameters $\mathbf{z} = g_{\theta}(\mathbf{z}_0)$. In EPI optimization, stochastic gradient steps
 329 in θ are taken such that entropy is maximized, and the emergent property \mathcal{X} is produced. The EPI posterior
 330 distribution is denoted $q_{\theta}(\mathbf{z} \mid \mathcal{X})$. **E.** The EPI posterior producing intermediate hub frequency. Samples are
 331 colored by log probability density. Distribution contours of average hub neuron frequency from mean of .55
 332 Hz are shown at levels of .525, .53,575 Hz (dark to light gray away from mean). Eigenvectors of the
 333 Hessian at the mode of the inferred distribution are indicated as \mathbf{v}_1 (solid) and \mathbf{v}_2 (dashed) with lengths
 334 scaled by the square root of the absolute value of their eigenvalues. **F** Simulations from parameters in E.
 335 (Top) The predictive distribution of the posterior obeys the emergent property. The black and gray dashed
 336 lines show the mean and two standard deviations according the emergent property, respectively. (Bottom)
 337 Simulations at the starred parameter values (figure 1F activity traces) (see Section 5.2.1 Stomatogastric
 338 ganglion subsection 5.2.1).

339 **3.2 A deep generative modeling approach to emergent property inference**

340 Emergent property inference (EPI) systematizes the three-step procedure of the previous section.
341 First, we consider the model as a coupled set of differential equations [57]. In the running STG
342 example, the model activity $\mathbf{x} = [x_{f1}, x_{f2}, x_{hub}, x_{s1}, x_{s2}]$ is the membrane potential for each neuron,
343 which evolves according to the biophysical conductance-based equation:

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

344 where $C_m = 1\text{nF}$, and \mathbf{h} is a sum of the leak, calcium, potassium, hyperpolarization, electrical, and
345 synaptic currents, all of which have their own complicated dependence on \mathbf{x} and $\mathbf{z} = [g_{el}, g_{synA}]$,
346 and $d\mathbf{B}$ is white gaussian noise (see Section 5.2.1Stomatogastric ganglionsubsubsection.5.2.1).

347 Second, we define the emergent property, which as above is “intermediate hub frequency” (Figure
348 1Emergent property inference (EPI) in the stomatogastric ganglion. **A.** Conductance-based biophysical
349 model of the STG subcircuit. In the STG model, jagged connections indicate electrical coupling having
350 electrical conductance g_{el} . Other connections in the diagram are inhibitory synaptic projections having
351 strength g_{synA} onto the hub neuron, and $g_{synB} = 5\text{nS}$ for mutual inhibitory connections. Parameters are
352 represented by the vector \mathbf{z} and membrane potentials by the vector \mathbf{x} . The evolution of this model’s activity
353 $\mathbf{x}(t)$ is predicated by differential equations. **B.** Spiking frequency $\omega(\mathbf{x}; \mathbf{z})$ is an emergent property statistic. In
354 this example, spiking frequency is measured from simulated activity of the STG model at parameter choices
355 of $g_{el} = 4.5\text{nS}$ and $g_{synA} = 3\text{nS}$. **C.** The emergent property of intermediate hub frequency, in which the hub
356 neuron fires at a rate between the fast and slow frequencies. This emergent property is defined by a mean and
357 variance on the emergent property statistic. Simulated activity traces are colored by log probability density
358 of their generating parameters in the EPI-inferred distribution (Panel E). **D.** For a choice of model and
359 emergent property, emergent property inference (EPI) learns a deep probability distribution of parameters
360 \mathbf{z} . Deep probability distributions map a simple random variable \mathbf{z}_0 through a deep neural network with
361 weights and biases $\boldsymbol{\theta}$ to parameters $\mathbf{z} = g_{\boldsymbol{\theta}}(\mathbf{z}_0)$. In EPI optimization, stochastic gradient steps in $\boldsymbol{\theta}$ are taken
362 such that entropy is maximized, and the emergent property \mathcal{X} is produced. The EPI posterior distribution is
363 denoted $q_{\boldsymbol{\theta}}(\mathbf{z} | \mathcal{X})$. **E.** The EPI posterior producing intermediate hub frequency. Samples are colored by log
364 probability density. Distribution contours of average hub neuron frequency from mean of .55 Hz are shown at
365 levels of .525, .53,575 Hz (dark to light gray away from mean). Eigenvectors of the Hessian at the mode
366 of the inferred distribution are indicated as \mathbf{v}_1 (solid) and \mathbf{v}_2 (dashed) with lengths scaled by the square
367 root of the absolute value of their eigenvalues. **F** Simulations from parameters in E. (Top) The predictive
368 distribution of the posterior obeys the emergent property. The black and gray dashed lines show the mean

369 and two standard deviations according the emergent property, respectively. (Bottom) Simulations at the
 370 starred parameter valuesfigure.1C). Quantifying this phenomenon is straightforward: we stipulate
 371 that the hub neuron’s spiking frequency – denoted $\omega_{\text{hub}}(\mathbf{x})$ is close to an intermediate frequency
 372 of 0.55Hz. Mathematically, we achieve this via constraints on the mean and variance of the hub
 373 neuron spiking frequency.

$$\begin{aligned} \mathcal{X} : \mathbb{E}_{\mathbf{z}, \mathbf{x}} [f(\mathbf{x}; \mathbf{z})] &\triangleq \mathbb{E}_{\mathbf{z}, \mathbf{x}} [\omega_{\text{hub}}(\mathbf{x}; \mathbf{z})] = [0.55] \triangleq \boldsymbol{\mu} \\ \text{Var}_{\mathbf{z}, \mathbf{x}} [f(\mathbf{x}; \mathbf{z})] &\triangleq \text{Var}_{\mathbf{z}, \mathbf{x}} [\omega_{\text{hub}}(\mathbf{x}; \mathbf{z})] = [0.025^2] \triangleq \boldsymbol{\sigma}^2. \end{aligned} \quad (2)$$

374 The emergent property statistic $f(\mathbf{x}; \mathbf{z}) = \omega_{\text{hub}}(\mathbf{x}; \mathbf{z})$ along with its constrained mean $\boldsymbol{\mu}$ and variance
 375 $\boldsymbol{\sigma}^2$ define the emergent property denoted \mathcal{X} .

376 Third, we perform emergent property inference: we find a distribution over parameter configura-
 377 tions \mathbf{z} , and insist that samples from this distribution produce the emergent property; in other
 378 words, they obey the constraints introduced in Equation 2A deep generative modeling approach to
 379 emergent property inferenceequation.3.2. This distribution will be chosen from a family of probabil-
 380 ity distributions $\mathcal{Q} = \{q_{\boldsymbol{\theta}}(\mathbf{z}) : \boldsymbol{\theta} \in \Theta\}$, defined by a deep generative distribution of the normalizing
 381 flow class [47, 50, 51] – neural networks which transform a simple distribution into a suitably com-
 382 plicated distribution (as is needed here). This deep distribution is represented in Figure 1Emergent
 383 property inference (EPI) in the stomatogastric ganglion. **A.** Conductance-based biophysical model of the
 384 STG subcircuit. In the STG model, jagged connections indicate electrical coupling having electrical con-
 385 ductance g_{el} . Other connections in the diagram are inhibitory synaptic projections having strength g_{synA}
 386 onto the hub neuron, and $g_{\text{synB}} = 5\text{nS}$ for mutual inhibitory connections. Parameters are represented by the
 387 vector \mathbf{z} and membrane potentials by the vector \mathbf{x} . The evolution of this model’s activity $\mathbf{x}(t)$ is predicated
 388 by differential equations. **B.** Spiking frequency $\omega(\mathbf{x}; \mathbf{z})$ is an emergent property statistic. In this example,
 389 spiking frequency is measured from simulated activity of the STG model at parameter choices of $g_{\text{el}} = 4.5\text{nS}$
 390 and $g_{\text{synA}} = 3\text{nS}$. **C.** The emergent property of intermediate hub frequency, in which the hub neuron fires
 391 at a rate between the fast and slow frequencies. This emergent property is defined by a mean and variance
 392 on the emergent property statistic. Simulated activity traces are colored by log probability density of their
 393 generating parameters in the EPI-inferred distribution (Panel E). **D.** For a choice of model and emergent
 394 property, emergent property inference (EPI) learns a deep probability distribution of parameters \mathbf{z} . Deep
 395 probability distributions map a simple random variable \mathbf{z}_0 through a deep neural network with weights and
 396 biases $\boldsymbol{\theta}$ to parameters $\mathbf{z} = g_{\boldsymbol{\theta}}(\mathbf{z}_0)$. In EPI optimization, stochastic gradient steps in $\boldsymbol{\theta}$ are taken such that
 397 entropy is maximized, and the emergent property \mathcal{X} is produced. The EPI posterior distribution is denoted

398 $q_{\theta}(\mathbf{z} \mid \mathcal{X})$. **E.** The EPI posterior producing intermediate hub frequency. Samples are colored by log prob-
 399 ability density. Distribution contours of average hub neuron frequency from mean of .55 Hz are shown at
 400 levels of .525, .53,575 Hz (dark to light gray away from mean). Eigenvectors of the Hessian at the mode
 401 of the inferred distribution are indicated as \mathbf{v}_1 (solid) and \mathbf{v}_2 (dashed) with lengths scaled by the square
 402 root of the absolute value of their eigenvalues. **F** Simulations from parameters in E. (Top) The predictive
 403 distribution of the posterior obeys the emergent property. The black and gray dashed lines show the mean
 404 and two standard deviations according the emergent property, respectively. (Bottom) Simulations at the
 405 starred parameter valuesfigure.1C (see Section 5.1Emergent property inference (EPI)subsection.5.1).
 406 Then, mathematically, we must solve the following optimization program:

$$\begin{aligned}
 q_{\theta}(\mathbf{z} \mid \mathcal{X}) &= \underset{\mathbf{q}_{\theta} \in \mathcal{Q}}{\operatorname{argmax}} H(q_{\theta}(\mathbf{z})) \\
 \text{s.t. } \mathcal{X} : \quad &\mathbb{E}_{\mathbf{z}, \mathbf{x}} [f(\mathbf{x}; \mathbf{z})] = \boldsymbol{\mu}, \operatorname{Var}_{\mathbf{z}, \mathbf{x}} [f(\mathbf{x}; \mathbf{z})] = \boldsymbol{\sigma}^2
 \end{aligned} \tag{3}$$

407 where $f(\mathbf{x}, \mathbf{z})$, $\boldsymbol{\mu}$, and $\boldsymbol{\sigma}$ are defined as in Equation ???. According to the emergent property of
 408 interest, $f(\mathbf{x}, \mathbf{z})$ may contain multiple statistics, in which case the mean and variance vectors $\boldsymbol{\mu}$
 409 and $\boldsymbol{\sigma}^2$ match this dimension. Finally, we recognize that many distributions in \mathcal{Q} will respect
 410 the emergent property constraints, so we select that which has maximum entropy. This principle,
 411 captured in Equation ?? by the primal objective H , identifies parameter distributions with minimal
 412 assumptions beyond some chosen structure [58, 59, 52, 60]. Such a normative principle of maximum
 413 entropy, which is also that of Bayesian inference, naturally fits with our scientific objective of
 414 reasoning about parametric sensitivity and robustness. The recovered distribution of EPI is as
 415 variable as possible along each parametric manifold such that it produces the emergent property.
 416 EPI optimizes the weights and biases $\boldsymbol{\theta}$ of the deep network (which induces the probability dis-
 417 tribution) by iteratively solving Equation ???. The optimization is complete when the sampled
 418 models with parameters $\mathbf{z} \sim q_{\theta}(z \mid \mathcal{X})$ produce activity consistent with the specified emergent
 419 property (Fig. S4). Such convergence is evaluated with a hypothesis test that the means and
 420 variances of each emergent property statistic are not different than their constrained values (see
 421 Section 5.1.3Augmented Lagrangian optimizationsubsubsection.5.1.3). Further validation of EPI is
 422 available in the supplementary materials, where we analyze a simpler model for which ground-truth
 423 statements can be made (Section 5.1.4Example: 2D LDSsubsubsection.5.1.4).

424 In relation to broader methodology, inspection of the EPI objective reveals a natural relationship
 425 to posterior inference. Specifically, EPI executes a novel variant of Bayesian inference with a uni-
 426 form prior and a gaussian likelihood on the emergent property statistic (see Section 5.1.5Maximum
 427 entropy distributions and exponential familiessubsection.5.1.5). A key advantage of EPI over

428 established Bayesian inference is that the predictions made by the inferred distribution are con-
429 strained to produce the specified emergent property. Equipped with this method, we may examine
430 structure in posterior distributions or make comparisons between posteriors conditioned at different
431 levels of the same emergent property statistic. In Sections 3.3EPI reveals how neuron-type specific
432 noise governs variability in the stochastic stabilized supralinear networkssubsection.3.3 and 3.4EPI
433 identifies multiple regimes of rapid task switchingsubsection.3.4, we prove out the value of EPI
434 by using it to investigate and produce novel insights into two prominent models in neuroscience.
435 Subsequently in Section 3.5EPI scales well to high-dimensional parameter spacessubsection.3.5, we
436 show EPI’s superiority in parameter scalability and fidelity of the posterior predictive distribution
437 by conditioning on stable amplification in low-rank RNNs.

438 **3.3 EPI reveals how neuron-type specific noise governs variability in the stochas-**
439 **tic stabilized supralinear network**

440 Dynamical models of excitatory (E) and inhibitory (I) populations with supralinear input-output
441 function have succeeded in explaining a host of experimentally documented phenomena. In a regime
442 characterized by inhibitory stabilization of strong recurrent excitation, these models give rise to
443 paradoxical responses [9], selective amplification [61, 62], surround suppression [63] and normal-
444 ization [64]. Despite their strong predictive power, E-I circuit models rely on the assumption that
445 inhibition can be studied as an indivisible unit. However, experimental evidence shows that inhibi-
446 tion is composed of distinct elements – parvalbumin (P), somatostatin (S), VIP (V) – composing
447 80% of GABAergic interneurons in V1 [65, 66, 67], and that these inhibitory cell types follow spe-
448 cific connectivity patterns (Fig. 2Emergent property inference in the stochastic stabilized supralinear
449 network (SSSN) **A**. Four-population model of primary visual cortex with excitatory (black), parvalbumin
450 (blue), somatostatin (red), and VIP (green) neurons (excitatory and inhibitory projections filled and unfilled,
451 respectively). Some neuron-types largely do not form synaptic projections to others ($|W_{\alpha_1, \alpha_2}| < 0.025$).
452 Each neural population receives a baseline input \mathbf{h}_b , and the E- and P-populations also receive a contrast-
453 dependent input \mathbf{h}_c . Additionally, each neural population receives a slow noisy input ϵ . **B**. Steady-state
454 responses of the SSN model (deterministic, $\sigma = 0$) to varying contrasts. The response at 50% contrast (dots)
455 is the focus of our analysis. **C**. Transient network responses of the SSSN model at 50 % contrast. (Left)
456 Traces are independent trials with varying initialization $\mathbf{x}(0)$ and noise realization. (Right) Mean (solid
457 line) and standard deviation (shading) of responses. **D**. EPI posterior of noise parameters \mathbf{z} conditioned
458 on E-population variability. The posterior predictive distribution of $s_E(\mathbf{x}; \mathbf{z})$ is show on the bottom-left.

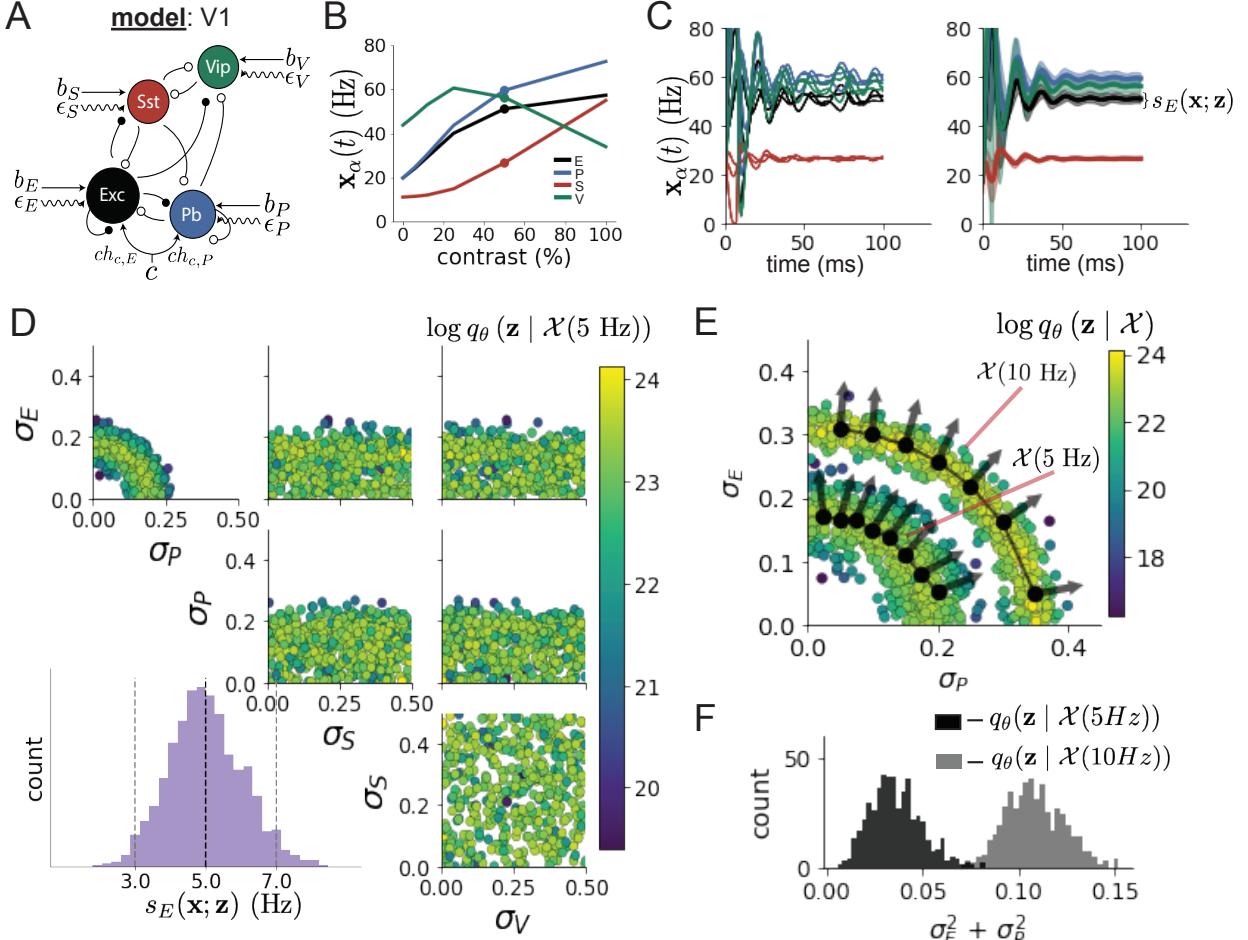


Figure 2: Emergent property inference in the stochastic stabilized supralinear network (SSSN) **A.** Four-population model of primary visual cortex with excitatory (black), parvalbumin (blue), somatostatin (red), and VIP (green) neurons (excitatory and inhibitory projections filled and unfilled, respectively). 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 ϵ . **B.** Steady-state responses of the SSN model (deterministic, $\sigma = \mathbf{0}$) to varying contrasts. The response at 50% contrast (dots) is the focus of our analysis. **C.** Transient network responses of the SSSN model at 50 % contrast. (Left) Traces are independent trials with varying initialization $\mathbf{x}(0)$ and noise realization. (Right) Mean (solid line) and standard deviation (shading) of responses. **D.** EPI posterior of noise parameters \mathbf{z} conditioned on E-population variability. The posterior predictive distribution of $s_E(\mathbf{x}; \mathbf{z})$ is show on the bottom-left. **E.** (Top) Enlarged visualization of the σ_E - σ_P marginal distribution of the posteriors $q_\theta(\mathbf{z} | \mathcal{X}(5 \text{ Hz}))$ and $q_\theta(\mathbf{z} | \mathcal{X}(10 \text{ Hz}))$. Each black dot shows the mode at each σ_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 posterior $q_\theta(\mathbf{z} | \mathcal{X}(5 \text{ Hz}))$ and $q_\theta(\mathbf{z} | \mathcal{X}(10 \text{ Hz}))$.

459 E. (Top) Enlarged visualization of the σ_E - σ_P marginal distribution of the posteriors $q_{\theta}(\mathbf{z} \mid \mathcal{X}(5 \text{ Hz})$ and
 460 $q_{\theta}(\mathbf{z} \mid \mathcal{X}(10 \text{ Hz})$. Each black dot shows the mode at each σ_P . The arrows show the most sensitive dimen-
 461 sions of the Hessian evaluated at these modes. F. The predictive distributions of $\sigma_E^2 + \sigma_P^2$ of each posterior
 462 $q_{\theta}(\mathbf{z} \mid \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} \mid \mathcal{X}(10 \text{ Hz}))$ [68]. Recent theoretical advances [54, 69, 70], have only
 463 started to address the consequences of this multiplicity in the dynamics of V1, strongly relying on
 464 linear theoretical tools. Here, we use EPI to analyze V1 models of greater complexity in order to
 465 characterize properties of slow noise governing circuit variability.

466 We considered the response properties of a nonlinear dynamical V1 circuit model (Fig. 2Emergent
 467 property inference in the stochastic stabilized supralinear network (SSSN) A. Four-population model of
 468 primary visual cortex with excitatory (black), parvalbumin (blue), somatostatin (red), and VIP (green)
 469 neurons (excitatory and inhibitory projections filled and unfilled, respectively). Some neuron-types largely
 470 do not form synaptic projections to others ($|W_{\alpha_1, \alpha_2}| < 0.025$). Each neural population receives a baseline
 471 input \mathbf{h}_b , and the E- and P-populations also receive a contrast-dependent input \mathbf{h}_c . Additionally, each neural
 472 population receives a slow noisy input ϵ . B. Steady-state responses of the SSN model (deterministic, $\sigma = 0$)
 473 to varying contrasts. The response at 50% contrast (dots) is the focus of our analysis. C. Transient network
 474 responses of the SSSN model at 50 % contrast. (Left) Traces are independent trials with varying initialization
 475 $\mathbf{x}(0)$ and noise realization. (Right) Mean (solid line) and standard deviation (shading) of responses. D. EPI
 476 posterior of noise parameters \mathbf{z} conditioned on E-population variability. The posterior predictive distribution
 477 of $s_E(\mathbf{x}; \mathbf{z})$ is show on the bottom-left. E. (Top) Enlarged visualization of the σ_E - σ_P marginal distribution of
 478 the posteriors $q_{\theta}(\mathbf{z} \mid \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} \mid \mathcal{X}(10 \text{ Hz}))$. Each black dot shows the mode at each σ_P . The arrows
 479 show the most sensitive dimensions of the Hessian evaluated at these modes. F. The predictive distributions
 480 of $\sigma_E^2 + \sigma_P^2$ of each posterior $q_{\theta}(\mathbf{z} \mid \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} \mid \mathcal{X}(10 \text{ Hz}))$ [68] with a state comprised of each
 481 neuron-type population's rate $\mathbf{x} = [x_E, x_P, x_S, x_V]^T$. Each population receives recurrent input $W\mathbf{x}$
 482 from synaptic projections of effective connectivity W and an external input \mathbf{h} , which determine
 483 the population rate via supralinear nonlinearity $\phi = \|\mathbf{x}\|_+^2$. The input is also comprised of a slow
 484 noise component $\epsilon \sim OU(\tau_{\text{noise}}, \sigma)$ of time scale $\tau_{\text{noise}} > \tau$ and variance parameters σ (see Section
 485 5.2.2Primary visual cortexsubsection.5.2.2)

$$\tau \frac{d\mathbf{x}}{dt} = -\mathbf{x} + \phi(W\mathbf{x} + \mathbf{h} + \epsilon). \quad (4)$$

486 This model is the stochastic stabilized supralinear network (SSSN) [71] generalized to have in-
 487 hibitory multiplicity, and introduces stochasticity to previous four neuron-type models of V1 [54].
 488 Stochasticity and inhibitory multiplicity introduce substantial complexity to mathematical deriva-
 489 tions (see Section 5.2.3Primary visual cortex: challenges to analysissubsection.5.2.3) motivating

490 the treatment of this model with EPI. Here, we consider fixed weights W and input \mathbf{h} according
 491 to a fit of the deterministic model to contrast responses [72] (Fig. 2Emergent property inference in
 492 the stochastic stabilized supralinear network (SSSN) **A**. Four-population model of primary visual cortex
 493 with excitatory (black), parvalbumin (blue), somatostatin (red), and VIP (green) neurons (excitatory and
 494 inhibitory projections filled and unfilled, respectively). Some neuron-types largely do not form synaptic pro-
 495 jections to others ($|W_{\alpha_1, \alpha_2}| < 0.025$). Each neural population receives a baseline input \mathbf{h}_b , and the E- and
 496 P-populations also receive a contrast-dependent input \mathbf{h}_c . Additionally, each neural population receives a
 497 slow noisy input ϵ . **B**. Steady-state responses of the SSN model (deterministic, $\sigma = \mathbf{0}$) to varying contrasts.
 498 The response at 50% contrast (dots) is the focus of our analysis. **C**. Transient network responses of the
 499 SSSN model at 50 % contrast. (Left) Traces are independent trials with varying initialization $\mathbf{x}(0)$ and noise
 500 realization. (Right) Mean (solid line) and standard deviation (shading) of responses. **D**. EPI posterior of
 501 noise parameters \mathbf{z} conditioned on E-population variability. The posterior predictive distribution of $s_E(\mathbf{x}; \mathbf{z})$
 502 is show on the bottom-left. **E**. (Top) Enlarged visualization of the σ_E - σ_P marginal distribution of the pos-
 503 teriors $q_{\theta}(\mathbf{z} | \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} | \mathcal{X}(10 \text{ Hz})$. Each black dot shows the mode at each σ_P . The arrows show
 504 the most sensitive dimensions of the Hessian evaluated at these modes. **F**. The predictive distributions of
 505 $\sigma_E^2 + \sigma_P^2$ of each posterior $q_{\theta}(\mathbf{z} | \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} | \mathcal{X}(10 \text{ Hz}))$ figure.2B), and study the effect of noise
 506 parameterization $\mathbf{z} = [\sigma_E, \sigma_P, \sigma_S, \sigma_V]^{\top}$ on fluctuations at 50% contrast.
 507 For this SSSN, we are interested in how noise variability across neural populations governs stochastic
 508 fluctuations in the E-population. Here, we quantify different levels y of E-population variability
 509 with the emergent property

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

510 where $s_E(\mathbf{x}; \mathbf{z})$ is the standard deviation of the stochastic E-population response about its steady
 511 state (Fig. 2Emergent property inference in the stochastic stabilized supralinear network (SSSN) **A**. Four-
 512 population model of primary visual cortex with excitatory (black), parvalbumin (blue), somatostatin (red),
 513 and VIP (green) neurons (excitatory and inhibitory projections filled and unfilled, respectively). Some
 514 neuron-types largely do not form synaptic projections to others ($|W_{\alpha_1, \alpha_2}| < 0.025$). Each neural population
 515 receives a baseline input \mathbf{h}_b , and the E- and P-populations also receive a contrast-dependent input \mathbf{h}_c .
 516 Additionally, each neural population receives a slow noisy input ϵ . **B**. Steady-state responses of the SSN
 517 model (deterministic, $\sigma = \mathbf{0}$) to varying contrasts. The response at 50% contrast (dots) is the focus of our
 518 analysis. **C**. Transient network responses of the SSSN model at 50 % contrast. (Left) Traces are independent
 519 trials with varying initialization $\mathbf{x}(0)$ and noise realization. (Right) Mean (solid line) and standard deviation

520 (shading) of responses. **D.** EPI posterior of noise parameters \mathbf{z} conditioned on E-population variability. The
 521 posterior predictive distribution of $s_E(\mathbf{x}; \mathbf{z})$ is show on the bottom-left. **E.** (Top) Enlarged visualization of the
 522 $\sigma_E\text{-}\sigma_P$ marginal distribution of the posteriors $q_{\theta}(\mathbf{z} | \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} | \mathcal{X}(10 \text{ Hz})$. Each black dot shows the
 523 mode at each σ_P . The arrows show the most sensitive dimensions of the Hessian evaluated at these modes.
 524 **F.** The predictive distributions of $\sigma_E^2 + \sigma_P^2$ of each posterior $q_{\theta}(\mathbf{z} | \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} | \mathcal{X}(10 \text{ Hz}))$ figure.2C).
 525 We ran EPI to obtain a posterior distribution $q_{\theta}(\mathbf{z} | \mathcal{X}(5 \text{ Hz})$ producing E-population variability
 526 around 5 Hz (Fig. 2Emergent property inference in the stochastic stabilized supralinear network (SSSN) **A.**
 527 Four-population model of primary visual cortex with excitatory (black), parvalbumin (blue), somatostatin
 528 (red), and VIP (green) neurons (excitatory and inhibitory projections filled and unfilled, respectively). Some
 529 neuron-types largely do not form synaptic projections to others ($|W_{\alpha_1, \alpha_2}| < 0.025$). Each neural population
 530 receives a baseline input \mathbf{h}_b , and the E- and P-populations also receive a contrast-dependent input \mathbf{h}_c .
 531 Additionally, each neural population receives a slow noisy input ϵ . **B.** Steady-state responses of the SSN
 532 model (deterministic, $\sigma = \mathbf{0}$) to varying contrasts. The response at 50% contrast (dots) is the focus of our
 533 analysis. **C.** Transient network responses of the SSSN model at 50 % contrast. (Left) Traces are independent
 534 trials with varying initialization $\mathbf{x}(0)$ and noise realization. (Right) Mean (solid line) and standard deviation
 535 (shading) of responses. **D.** EPI posterior of noise parameters \mathbf{z} conditioned on E-population variability. The
 536 posterior predictive distribution of $s_E(\mathbf{x}; \mathbf{z})$ is show on the bottom-left. **E.** (Top) Enlarged visualization
 537 of the $\sigma_E\text{-}\sigma_P$ marginal distribution of the posteriors $q_{\theta}(\mathbf{z} | \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} | \mathcal{X}(10 \text{ Hz})$. Each black
 538 dot shows the mode at each σ_P . The arrows show the most sensitive dimensions of the Hessian evaluated
 539 at these modes. **F.** The predictive distributions of $\sigma_E^2 + \sigma_P^2$ of each posterior $q_{\theta}(\mathbf{z} | \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} |$
 540 $\mathcal{X}(10 \text{ Hz}))$ figure.2D). From the marginal distribution of σ_E and σ_P (Fig. 2Emergent property inference
 541 in the stochastic stabilized supralinear network (SSSN) **A.** Four-population model of primary visual cortex
 542 with excitatory (black), parvalbumin (blue), somatostatin (red), and VIP (green) neurons (excitatory and
 543 inhibitory projections filled and unfilled, respectively). Some neuron-types largely do not form synaptic
 544 projections to others ($|W_{\alpha_1, \alpha_2}| < 0.025$). Each neural population receives a baseline input \mathbf{h}_b , and the E-
 545 and P-populations also receive a contrast-dependent input \mathbf{h}_c . Additionally, each neural population receives a
 546 slow noisy input ϵ . **B.** Steady-state responses of the SSN model (deterministic, $\sigma = \mathbf{0}$) to varying contrasts.
 547 The response at 50% contrast (dots) is the focus of our analysis. **C.** Transient network responses of the
 548 SSSN model at 50 % contrast. (Left) Traces are independent trials with varying initialization $\mathbf{x}(0)$ and noise
 549 realization. (Right) Mean (solid line) and standard deviation (shading) of responses. **D.** EPI posterior of
 550 noise parameters \mathbf{z} conditioned on E-population variability. The posterior predictive distribution of $s_E(\mathbf{x}; \mathbf{z})$
 551 is show on the bottom-left. **E.** (Top) Enlarged visualization of the $\sigma_E\text{-}\sigma_P$ marginal distribution of the

552 posteriors $q_{\theta}(\mathbf{z} \mid \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} \mid \mathcal{X}(10 \text{ Hz})$. Each black dot shows the mode at each σ_P . The arrows
 553 show the most sensitive dimensions of the Hessian evaluated at these modes. **F.** The predictive distributions
 554 of $\sigma_E^2 + \sigma_P^2$ of each posterior $q_{\theta}(\mathbf{z} \mid \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} \mid \mathcal{X}(10 \text{ Hz}))$ figure.2D, top-left), we can see that
 555 $s_E(\mathbf{x}; \mathbf{z})$ is sensitive to various combinations of σ_E and σ_P . Alternatively, both σ_S and σ_V are
 556 degenerate with respect to $s_E(\mathbf{x}; \mathbf{z})$ evidenced by the high variability in those dimensions of the
 557 posterior (Fig. 2Emergent property inference in the stochastic stabilized supralinear network (SSSN) **A**.
 558 Four-population model of primary visual cortex with excitatory (black), parvalbumin (blue), somatostatin
 559 (red), and VIP (green) neurons (excitatory and inhibitory projections filled and unfilled, respectively). Some
 560 neuron-types largely do not form synaptic projections to others ($|W_{\alpha_1, \alpha_2}| < 0.025$). Each neural population
 561 receives a baseline input \mathbf{h}_b , and the E- and P-populations also receive a contrast-dependent input \mathbf{h}_c .
 562 Additionally, each neural population receives a slow noisy input ϵ . **B.** Steady-state responses of the SSN
 563 model (deterministic, $\sigma = \mathbf{0}$) to varying contrasts. The response at 50% contrast (dots) is the focus of our
 564 analysis. **C.** Transient network responses of the SSSN model at 50 % contrast. (Left) Traces are independent
 565 trials with varying initialization $\mathbf{x}(0)$ and noise realization. (Right) Mean (solid line) and standard deviation
 566 (shading) of responses. **D.** EPI posterior of noise parameters \mathbf{z} conditioned on E-population variability. The
 567 posterior predictive distribution of $s_E(\mathbf{x}; \mathbf{z})$ is show on the bottom-left. **E.** (Top) Enlarged visualization of the
 568 σ_E - σ_P marginal distribution of the posteriors $q_{\theta}(\mathbf{z} \mid \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} \mid \mathcal{X}(10 \text{ Hz}))$. Each black dot shows the
 569 mode at each σ_P . The arrows show the most sensitive dimensions of the Hessian evaluated at these modes.
 570 **F.** The predictive distributions of $\sigma_E^2 + \sigma_P^2$ of each posterior $q_{\theta}(\mathbf{z} \mid \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} \mid \mathcal{X}(10 \text{ Hz}))$ figure.2D,
 571 bottom-right). Together, these observations imply a parametric manifold of degeneracy with respect
 572 to $s_E(\mathbf{x}; \mathbf{z})$ of 5 Hz, which is indicated by the modes along σ_P in the σ_E - σ_P marginal (Fig. 2Emergent
 573 property inference in the stochastic stabilized supralinear network (SSSN) **A**. Four-population model of
 574 primary visual cortex with excitatory (black), parvalbumin (blue), somatostatin (red), and VIP (green)
 575 neurons (excitatory and inhibitory projections filled and unfilled, respectively). Some neuron-types largely
 576 do not form synaptic projections to others ($|W_{\alpha_1, \alpha_2}| < 0.025$). Each neural population receives a baseline
 577 input \mathbf{h}_b , and the E- and P-populations also receive a contrast-dependent input \mathbf{h}_c . Additionally, each
 578 neural population receives a slow noisy input ϵ . **B.** Steady-state responses of the SSN model (deterministic,
 579 $\sigma = \mathbf{0}$) to varying contrasts. The response at 50% contrast (dots) is the focus of our analysis. **C.** Transient
 580 network responses of the SSSN model at 50 % contrast. (Left) Traces are independent trials with varying
 581 initialization $\mathbf{x}(0)$ and noise realization. (Right) Mean (solid line) and standard deviation (shading) of
 582 responses. **D.** EPI posterior of noise parameters \mathbf{z} conditioned on E-population variability. The posterior
 583 predictive distribution of $s_E(\mathbf{x}; \mathbf{z})$ is show on the bottom-left. **E.** (Top) Enlarged visualization of the σ_E - σ_P
 584 marginal distribution of the posteriors $q_{\theta}(\mathbf{z} \mid \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} \mid \mathcal{X}(10 \text{ Hz}))$. Each black dot shows the mode

at each σ_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 posterior $q_{\theta}(\mathbf{z} | \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} | \mathcal{X}(10 \text{ Hz}))$ figure.2E).
 The dimensions of sensitivity conferred by EPI and this plain visual structure suggest a quadratic
 relationship in the emergent property statistic $s_E(\mathbf{x}; \mathbf{z})$ and parameters \mathbf{z} , which is preserved at
 a greater level of variability $\mathcal{X}(10 \text{ Hz})$ (Fig. 2Emergent 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 ϵ . **B.** Steady-state responses of the SSN model (deterministic, $\sigma = \mathbf{0}$) to varying contrasts. The
 response at 50% contrast (dots) is the focus of our analysis. **C.** Transient network responses of the SSSN
 model at 50 % contrast. (Left) Traces are independent trials with varying initialization $\mathbf{x}(0)$ and noise
 realization. (Right) Mean (solid line) and standard deviation (shading) of responses. **D.** EPI posterior of
 noise parameters \mathbf{z} conditioned on E-population variability. The posterior predictive distribution of $s_E(\mathbf{x}; \mathbf{z})$
 is show on the bottom-left. **E.** (Top) Enlarged visualization of the σ_E - σ_P marginal distribution of the
 posteriors $q_{\theta}(\mathbf{z} | \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} | \mathcal{X}(10 \text{ Hz}))$. Each black dot shows the mode at each σ_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 posterior $q_{\theta}(\mathbf{z} | \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} | \mathcal{X}(10 \text{ Hz}))$ figure.2E). Indeed, the sum of squares
 of σ_E and σ_P is larger in $q_{\theta}(\mathbf{z} | \mathcal{X}(10 \text{ Hz}))$ than $q_{\theta}(\mathbf{z} | \mathcal{X}(5 \text{ Hz}))$ (Fig 2Emergent 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 ϵ . **B.** Steady-state responses of the SSN model (deterministic, $\sigma = \mathbf{0}$) to varying contrasts.
 The response at 50% contrast (dots) is the focus of our analysis. **C.** Transient network responses of the
 SSSN model at 50 % contrast. (Left) Traces are independent trials with varying initialization $\mathbf{x}(0)$ and noise
 realization. (Right) Mean (solid line) and standard deviation (shading) of responses. **D.** EPI posterior of
 noise parameters \mathbf{z} conditioned on E-population variability. The posterior predictive distribution of $s_E(\mathbf{x}; \mathbf{z})$
 is show on the bottom-left. **E.** (Top) Enlarged visualization of the σ_E - σ_P marginal distribution of the
 posteriors $q_{\theta}(\mathbf{z} | \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} | \mathcal{X}(10 \text{ Hz}))$. Each black dot shows the mode at each σ_P . The arrows
 show the most sensitive dimensions of the Hessian evaluated at these modes. **F.** The predictive distributions

618 of $\sigma_E^2 + \sigma_P^2$ of each posterior $q_{\theta}(\mathbf{z} \mid \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} \mid \mathcal{X}(10 \text{ Hz}))$ figure.2F, $p = 0$), while the sum
619 of squares of σ_S and σ_V are not significantly different in the two posteriors (Fig. 11(V1 3) EPI
620 posterior for $\mathcal{X}(10 \text{ Hz})$ figure.11, $p = .402$).

621 While a quadratic relationship in $s_E(\mathbf{x}; \mathbf{z})$ and \mathbf{z} is potentially derivable by extending the derivation
622 in Section 5.2.2Primary visual cortexsubsubsection.5.2.2 to the case of $\tau \neq \tau_{\text{noise}}$, the coefficients
623 in front of each quadratic term would be unruly, and likely escape comprehensible analysis. This
624 makes EPI an attractive tool for revealing the characteristics of noise governing variability and for
625 answering other questions in this complex model. Intriguingly, this circuit exhibited a paradoxical
626 effect in the P-population, and no other inhibitory types at 50% contrast (Fig. 11(V1 3) EPI posterior
627 for $\mathcal{X}(10 \text{ Hz})$ figure.11) implying that the E-population is P-stabilized. Future work motivated by
628 our analysis here, may uncover a relationship between the neuron-type mediating stability and the
629 factors governing circuit variability.

630 3.4 EPI identifies multiple regimes of rapid task switching

631 In a rapid task switching experiment [73], rats were explicitly cued on each trial to either orient
632 towards a visual stimulus in the Pro (P) task or orient away from a visual stimulus in the Anti
633 (A) task (Fig. 3A. Rapid task switching behavioral paradigm (see text). B. Model of superior colliculus (SC).
634 Neurons: LP - left pro, RP - right pro, LA - left anti, RA - right anti. Parameters: sW - self, hW - horizontal, vW
635 -vertical, dW - diagonal weights. C. The EPI posterior distribution of rapid task switching networks. Red and purple
636 stars (\mathbf{z}_1 and \mathbf{z}_2) indicate different connectivity regimes with different sensitivity vectors \mathbf{v}_1 and \mathbf{v}_2 . (Middle-left)
637 Posterior predictive distribution of task accuracies. (Bottom-left) Task accuracy along dimensions of sensitivity in
638 each connectivity regime. D. Means (solid) and standard deviations (shaded) of each population across random
639 simulated trials. Top plots show Pro (top) and Anti (bottom) responses for connectivity \mathbf{z}_1 . Bottom rows show the
640 same \mathbf{z}_2 . E. The EPI posterior predicts experimental results (left) showing no change in the Pro task, but larger error
641 in the Anti task (right). F. Accuracy in the Anti task during delay period optogenetic inactivation $p_{A,\text{opto}}$ is strongly
642 anticorrelated with accuracy in the Pro task. G. Accuracy with delay period inactivation along each connectivity
643 regime's dimension of sensitivityfigure.3A). Neural recordings in the midbrain superior colliculus (SC)
644 exhibited two populations of neurons that simultaneously represented both task context (Pro or
645 Anti) and motor response (contralateral or ipsilateral to the recorded side): the Pro/Contra and
646 Anti/Ipsi neurons [55]. Duan et al. proposed a model of SC that, like the V1 model analyzed
647 in the previous section, is a four-population dynamical system. We analyzed this model, where
648 the neuron-type populations are functionally-defined as the Pro- and Anti-populations in each

649 hemisphere (left (L) and right (R)), their connectivity is parameterized geometrically (Fig. 3A.
 650 Rapid task switching behavioral paradigm (see text). **B.** Model of superior colliculus (SC). Neurons: LP - left pro,
 651 RP - right pro, LA - left anti, RA - right anti. Parameters: sW - self, hW - horizontal, vW -vertical, dW - diagonal
 652 weights. **C.** The EPI posterior distribution of rapid task switching networks. Red and purple stars (\mathbf{z}_1 and \mathbf{z}_2)
 653 indicate different connectivity regimes with different sensitivity vectors \mathbf{v}_1 and \mathbf{v}_2 . (Middle-left) Posterior predictive
 654 distribution of task accuracies. (Bottom-left) Task accuracy along dimensions of sensitivity in each connectivity
 655 regime. **D.** Means (solid) and standard deviations (shaded) of each population across random simulated trials. Top
 656 plots show Pro (top) and Anti (bottom) responses for connectivity \mathbf{z}_1 . Bottom rows show the same \mathbf{z}_2 . **E.** The EPI
 657 posterior predicts experimental results (left) showing no change in the Pro task, but larger error in the Anti task
 658 (right). **F.** Accuracy in the Anti task during delay period optogenetic inactivation $p_{A,\text{opto}}$ is strongly anticorrelated
 659 with accuracy in the Pro task. **G.** Accuracy with delay period inactivation along each connectivity regime's dimension
 660 of sensitivityfigure.3B). The input-output function of this model is chosen such that the population
 661 responses $\mathbf{x} = [x_{LP}, x_{LA}, x_{RP}, x_{RA}]^\top$ are bounded from 0 to 1 as a function ϕ of a dynamically
 662 evolving internal variable \mathbf{u} . The model responds to the side with greater Pro neuron activation;
 663 e.g. the response is left if $x_{LP} > x_{RP}$ at the end of the trial. The dynamics evolve with timescale
 664 $\tau = 90\text{ms}$ governed by connectivity weights W

$$\begin{aligned}\tau \frac{d\mathbf{u}}{dt} &= -\mathbf{u} + W\mathbf{x} + \mathbf{h} + d\mathbf{B} \\ \mathbf{x} &= \phi(\mathbf{u})\end{aligned}\tag{6}$$

665 with white noise of variance 0.2^2 . The input \mathbf{h} is comprised of a cue-dependent input to the Pro
 666 or Anti populations, a stimulus orientation input to either the Left or Right populations, and a
 667 choice-period input to the entire network (see Section 5.2.4Superior colliculussubsection.5.2.4).
 668 Here, we use EPI to determine the network connectivity $\mathbf{z} = [sW, vW, dW, hW]^\top$ that produces
 669 rapid task switching behavior.

670 We define rapid task switching behavior as accurate execution of each task. Inferred models should
 671 not exhibit fully random responses (50%), or perfect performance (100%), since perfection is never
 672 attained by even the best trained rats. We formulate rapid task switching as an emergent property
 673 by stipulating that the average accuracy in the Pro task $p_P(\mathbf{x}; \mathbf{z})$ and Anti task $p_A(\mathbf{x}; \mathbf{z})$ be 75%
 674 with variance $7.5\%^2$.

$$\begin{aligned}\mathcal{X} : \mathbb{E}_{\mathbf{z}} \begin{bmatrix} p_P(\mathbf{x}; \mathbf{z}) \\ p_A(\mathbf{x}; \mathbf{z}) \end{bmatrix} &= \begin{bmatrix} 75\% \\ 75\% \end{bmatrix} \\ \text{Var}_{\mathbf{z}} \begin{bmatrix} p_P(\mathbf{x}; \mathbf{z}) \\ p_A(\mathbf{x}; \mathbf{z}) \end{bmatrix} &= \begin{bmatrix} 7.5\%^2 \\ 7.5\%^2 \end{bmatrix}\end{aligned}\tag{7}$$

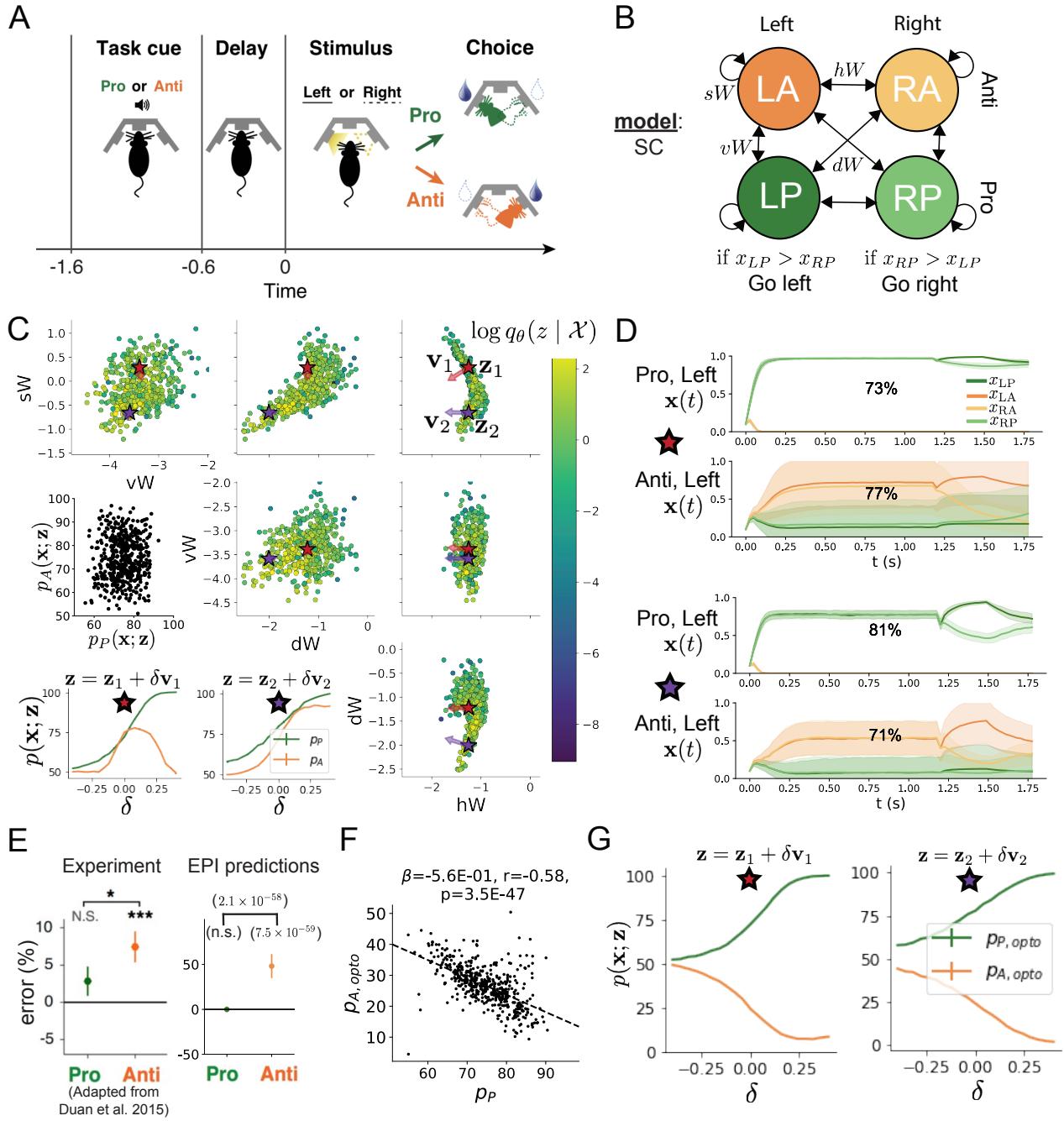


Figure 3: **A.** Rapid task switching behavioral paradigm (see text). **B.** Model of superior colliculus (SC). Neurons: LP - left pro, RP - right pro, LA - left anti, RA - right anti. Parameters: sW - self, hW - horizontal, vW - vertical, dW - diagonal weights. **C.** The EPI posterior distribution of rapid task switching networks. Red and purple stars (\mathbf{z}_1 and \mathbf{z}_2) indicate different connectivity regimes with different sensitivity vectors \mathbf{v}_1 and \mathbf{v}_2 . (Middle-left) Posterior predictive distribution of task accuracies. (Bottom-left) Task accuracy along dimensions of sensitivity in each connectivity regime. **D.** Means (solid) and standard deviations (shaded) of each population across random simulated trials. Top plots show Pro (top) and Anti (bottom) responses for connectivity \mathbf{z}_1 . Bottom rows show the same \mathbf{z}_2 . **E.** The EPI posterior predicts experimental results (left) showing no change in the Pro task, but larger error in the Anti task (right). **F.** Accuracy in the Anti task during delay period optogenetic inactivation $p_{A, \text{opto}}$ is strongly anticorrelated with accuracy in the Pro task. **G.** Accuracy with delay period inactivation along each connectivity regime's dimension of sensitivity.

675 A variance of 7.5%² in each task will confer a posterior producing performances ranging from about
676 60% – 90%, allowing us to examine the properties of connectivity that yield better performance in
677 each task. Notably, this is our first example using EPI to condition on multiple emergent property
678 statistics ($|f(\mathbf{x}; \mathbf{z})| = 2$).

679 The EPI inferred parameters (Fig. 3**A**. Rapid task switching behavioral paradigm (see text). **B**. Model
680 of superior colliculus (SC). Neurons: LP - left pro, RP - right pro, LA - left anti, RA - right anti. Parameters:
681 sW - self, hW - horizontal, vW - vertical, dW - diagonal weights. **C**. The EPI posterior distribution of rapid task
682 switching networks. Red and purple stars (\mathbf{z}_1 and \mathbf{z}_2) indicate different connectivity regimes with different sensitivity
683 vectors \mathbf{v}_1 and \mathbf{v}_2 . (Middle-left) Posterior predictive distribution of task accuracies. (Bottom-left) Task accuracy
684 along dimensions of sensitivity in each connectivity regime. **D**. Means (solid) and standard deviations (shaded)
685 of each population across random simulated trials. Top plots show Pro (top) and Anti (bottom) responses for
686 connectivity \mathbf{z}_1 . Bottom rows show the same \mathbf{z}_2 . **E**. The EPI posterior predicts experimental results (left) showing
687 no change in the Pro task, but larger error in the Anti task (right). **F**. Accuracy in the Anti task during delay period
688 optogenetic inactivation $p_{A,\text{opto}}$ is strongly anticorrelated with accuracy in the Pro task. **G**. Accuracy with delay
689 period inactivation along each connectivity regime's dimension of sensitivityfigure.3C) generate a distribution
690 of task accuracies (Fig. 3**A**. Rapid task switching behavioral paradigm (see text). **B**. Model of superior colliculus
691 (SC). Neurons: LP - left pro, RP - right pro, LA - left anti, RA - right anti. Parameters: sW - self, hW - horizontal,
692 vW - vertical, dW - diagonal weights. **C**. The EPI posterior distribution of rapid task switching networks. Red and
693 purple stars (\mathbf{z}_1 and \mathbf{z}_2) indicate different connectivity regimes with different sensitivity vectors \mathbf{v}_1 and \mathbf{v}_2 . (Middle-
694 left) Posterior predictive distribution of task accuracies. (Bottom-left) Task accuracy along dimensions of sensitivity
695 in each connectivity regime. **D**. Means (solid) and standard deviations (shaded) of each population across random
696 simulated trials. Top plots show Pro (top) and Anti (bottom) responses for connectivity \mathbf{z}_1 . Bottom rows show the
697 same \mathbf{z}_2 . **E**. The EPI posterior predicts experimental results (left) showing no change in the Pro task, but larger error
698 in the Anti task (right). **F**. Accuracy in the Anti task during delay period optogenetic inactivation $p_{A,\text{opto}}$ is strongly
699 anticorrelated with accuracy in the Pro task. **G**. Accuracy with delay period inactivation along each connectivity
700 regime's dimension of sensitivityfigure.3C, middle-left) according to our mathematical definition of rapid
701 task switching (Equation 4EPI identifies multiple regimes of rapid task switchingequation.3.4). The
702 nonlinear patterns of connectivity that govern each task accuracy (Fig. 12(SC1): **A**. Same pairplot
703 as Fig. 3**A**. Rapid task switching behavioral paradigm (see text). **B**. Model of superior colliculus (SC). Neurons:
704 LP - left pro, RP - right pro, LA - left anti, RA - right anti. Parameters: sW - self, hW - horizontal, vW -
705 vertical, dW - diagonal weights. **C**. The EPI posterior distribution of rapid task switching networks. Red and purple
706 stars (\mathbf{z}_1 and \mathbf{z}_2) indicate different connectivity regimes with different sensitivity vectors \mathbf{v}_1 and \mathbf{v}_2 . (Middle-left)

707 Posterior predictive distribution of task accuracies. (Bottom-left) Task accuracy along dimensions of sensitivity in
 708 each connectivity regime. **D.** Means (solid) and standard deviations (shaded) of each population across random
 709 simulated trials. Top plots show Pro (top) and Anti (bottom) responses for connectivity \mathbf{z}_1 . Bottom rows show the
 710 same \mathbf{z}_2 . **E.** The EPI posterior predicts experimental results (left) showing no change in the Pro task, but larger error
 711 in the Anti task (right). **F.** Accuracy in the Anti task during delay period optogenetic inactivation $p_{A,\text{opto}}$ is strongly
 712 anticorrelated with accuracy in the Pro task. **G.** Accuracy with delay period inactivation along each connectivity
 713 regime's dimension of sensitivity figure.3C colored by Pro task accuracy. **B.** Same as A colored by Anti task
 714 accuracy. **C.** Connectivity parameters of EPI distributions versus task accuracies. β is slope coefficient of
 715 linear regression, r is correlation, and p is the two-tailed p-value figure.12A-B) are not fully captured by
 716 linear prediction (Fig. 12(SC1): **A.** Same pairplot as Fig. 3A. Rapid task switching behavioral paradigm (see
 717 text). **B.** Model of superior colliculus (SC). Neurons: LP - left pro, RP - right pro, LA - left anti, RA - right anti.
 718 Parameters: sW - self, hW - horizontal, vW - vertical, dW - diagonal weights. **C.** The EPI posterior distribution
 719 of rapid task switching networks. Red and purple stars (\mathbf{z}_1 and \mathbf{z}_2) indicate different connectivity regimes with
 720 different sensitivity vectors \mathbf{v}_1 and \mathbf{v}_2 . (Middle-left) Posterior predictive distribution of task accuracies. (Bottom-
 721 left) Task accuracy along dimensions of sensitivity in each connectivity regime. **D.** Means (solid) and standard
 722 deviations (shaded) of each population across random simulated trials. Top plots show Pro (top) and Anti (bottom)
 723 responses for connectivity \mathbf{z}_1 . Bottom rows show the same \mathbf{z}_2 . **E.** The EPI posterior predicts experimental results
 724 (left) showing no change in the Pro task, but larger error in the Anti task (right). **F.** Accuracy in the Anti task
 725 during delay period optogenetic inactivation $p_{A,\text{opto}}$ is strongly anticorrelated with accuracy in the Pro task. **G.**
 726 Accuracy with delay period inactivation along each connectivity regime's dimension of sensitivity figure.3C colored
 727 by Pro task accuracy. **B.** Same as A colored by Anti task accuracy. **C.** Connectivity parameters of EPI
 728 distributions versus task accuracies. β is slope coefficient of linear regression, r is correlation, and p is the two-
 729 tailed p-value figure.12C). For example, the patterns in connectivity increasing Pro accuracy change
 730 dramatically after crossing a threshold of sW (Fig. 12(SC1): **A.** Same pairplot as Fig. 3A. Rapid task
 731 switching behavioral paradigm (see text). **B.** Model of superior colliculus (SC). Neurons: LP - left pro, RP - right
 732 pro, LA - left anti, RA - right anti. Parameters: sW - self, hW - horizontal, vW - vertical, dW - diagonal weights. **C.**
 733 The EPI posterior distribution of rapid task switching networks. Red and purple stars (\mathbf{z}_1 and \mathbf{z}_2) indicate different
 734 connectivity regimes with different sensitivity vectors \mathbf{v}_1 and \mathbf{v}_2 . (Middle-left) Posterior predictive distribution of
 735 task accuracies. (Bottom-left) Task accuracy along dimensions of sensitivity in each connectivity regime. **D.** Means
 736 (solid) and standard deviations (shaded) of each population across random simulated trials. Top plots show Pro (top)
 737 and Anti (bottom) responses for connectivity \mathbf{z}_1 . Bottom rows show the same \mathbf{z}_2 . **E.** The EPI posterior predicts
 738 experimental results (left) showing no change in the Pro task, but larger error in the Anti task (right). **F.** Accuracy in
 739 the Anti task during delay period optogenetic inactivation $p_{A,\text{opto}}$ is strongly anticorrelated with accuracy in the Pro

740 task. **G.** Accuracy with delay period inactivation along each connectivity regime's dimension of sensitivityfigure.3C
741 colored by Pro task accuracy. **B.** Same as A colored by Anti task accuracy. **C.** Connectivity parameters
742 of EPI distributions versus task accuracies. β is slope coefficient of linear regression, r is correlation, and
743 p is the two-tailed p-valuefigure.12A sW - hW marginal). Not only has EPI captured this complex
744 nonlinear posterior, it offers probabilistic tools for understanding the different regimes of model
745 behavior.

746 To establish these two regimes of connectivity, we took gradient steps along $q_{\theta}(\mathbf{z} \mid \mathcal{X})$ to produce
747 modes \mathbf{z}_1 and \mathbf{z}_2 (Fig. 3A. Rapid task switching behavioral paradigm (see text). **B.** Model of superior colliculus
748 (SC). Neurons: LP - left pro, RP - right pro, LA - left anti, RA - right anti. Parameters: sW - self, hW - hori-
749 zontal, vW -vertical, dW - diagonal weights. **C.** The EPI posterior distribution of rapid task switching networks.
750 Red and purple stars (\mathbf{z}_1 and \mathbf{z}_2) indicate different connectivity regimes with different sensitivity vectors \mathbf{v}_1 and \mathbf{v}_2 .
751 (Middle-left) Posterior predictive distribution of task accuracies. (Bottom-left) Task accuracy along dimensions of
752 sensitivity in each connectivity regime. **D.** Means (solid) and standard deviations (shaded) of each population across
753 random simulated trials. Top plots show Pro (top) and Anti (bottom) responses for connectivity \mathbf{z}_1 . Bottom rows
754 show the same \mathbf{z}_2 . **E.** The EPI posterior predicts experimental results (left) showing no change in the Pro task,
755 but larger error in the Anti task (right). **F.** Accuracy in the Anti task during delay period optogenetic inactivation
756 $p_{A,\text{opto}}$ is strongly anticorrelated with accuracy in the Pro task. **G.** Accuracy with delay period inactivation along
757 each connectivity regime's dimension of sensitivityfigure.3C red and purple stars, Section 5.2.4Superior
758 colliculussubsubsection.5.2.4). Simulations from these two regimes reveal different responses in
759 each task (Fig. 3A. Rapid task switching behavioral paradigm (see text). **B.** Model of superior colliculus (SC).
760 Neurons: LP - left pro, RP - right pro, LA - left anti, RA - right anti. Parameters: sW - self, hW - horizontal, vW
761 -vertical, dW - diagonal weights. **C.** The EPI posterior distribution of rapid task switching networks. Red and purple
762 stars (\mathbf{z}_1 and \mathbf{z}_2) indicate different connectivity regimes with different sensitivity vectors \mathbf{v}_1 and \mathbf{v}_2 . (Middle-left)
763 Posterior predictive distribution of task accuracies. (Bottom-left) Task accuracy along dimensions of sensitivity in
764 each connectivity regime. **D.** Means (solid) and standard deviations (shaded) of each population across random
765 simulated trials. Top plots show Pro (top) and Anti (bottom) responses for connectivity \mathbf{z}_1 . Bottom rows show the
766 same \mathbf{z}_2 . **E.** The EPI posterior predicts experimental results (left) showing no change in the Pro task, but larger error
767 in the Anti task (right). **F.** Accuracy in the Anti task during delay period optogenetic inactivation $p_{A,\text{opto}}$ is strongly
768 anticorrelated with accuracy in the Pro task. **G.** Accuracy with delay period inactivation along each connectivity
769 regime's dimension of sensitivityfigure.3D). We charcaterized these regimes by identifying the dimensions
770 of connectivity that rapid task switching is most sensitive to. The sensitivity dimensions \mathbf{v}_1 and \mathbf{v}_2
771 (Fig. 3A. Rapid task switching behavioral paradigm (see text). **B.** Model of superior colliculus (SC). Neurons: LP

772 - left pro, RP - right pro, LA - left anti, RA - right anti. Parameters: sW - self, hW - horizontal, vW -vertical, dW -
 773 diagonal weights. **C.** The EPI posterior distribution of rapid task switching networks. Red and purple stars (\mathbf{z}_1 and
 774 \mathbf{z}_2) indicate different connectivity regimes with different sensitivity vectors \mathbf{v}_1 and \mathbf{v}_2 . (Middle-left) Posterior predic-
 775 tive distribution of task accuracies. (Bottom-left) Task accuracy along dimensions of sensitivity in each connectivity
 776 regime. **D.** Means (solid) and standard deviations (shaded) of each population across random simulated trials. Top
 777 plots show Pro (top) and Anti (bottom) responses for connectivity \mathbf{z}_1 . Bottom rows show the same \mathbf{z}_2 . **E.** The EPI
 778 posterior predicts experimental results (left) showing no change in the Pro task, but larger error in the Anti task
 779 (right). **F.** Accuracy in the Anti task during delay period optogenetic inactivation $p_{A,\text{opto}}$ is strongly anticorrelated
 780 with accuracy in the Pro task. **G.** Accuracy with delay period inactivation along each connectivity regime's dimen-
 781 sion of sensitivityfigure.3C, red and purple arrows) point in different directions, resulting in different
 782 changes to task accuracy (Fig. 3**A**. Rapid task switching behavioral paradigm (see text). **B.** Model of superior
 783 colliculus (SC). Neurons: LP - left pro, RP - right pro, LA - left anti, RA - right anti. Parameters: sW - self, hW -
 784 horizontal, vW -vertical, dW - diagonal weights. **C.** The EPI posterior distribution of rapid task switching networks.
 785 Red and purple stars (\mathbf{z}_1 and \mathbf{z}_2) indicate different connectivity regimes with different sensitivity vectors \mathbf{v}_1 and \mathbf{v}_2 .
 786 (Middle-left) Posterior predictive distribution of task accuracies. (Bottom-left) Task accuracy along dimensions of
 787 sensitivity in each connectivity regime. **D.** Means (solid) and standard deviations (shaded) of each population across
 788 random simulated trials. Top plots show Pro (top) and Anti (bottom) responses for connectivity \mathbf{z}_1 . Bottom rows
 789 show the same \mathbf{z}_2 . **E.** The EPI posterior predicts experimental results (left) showing no change in the Pro task,
 790 but larger error in the Anti task (right). **F.** Accuracy in the Anti task during delay period optogenetic inactivation
 791 $p_{A,\text{opto}}$ is strongly anticorrelated with accuracy in the Pro task. **G.** Accuracy with delay period inactivation along
 792 each connectivity regime's dimension of sensitivityfigure.3D, bottom-left, 13(SC2): **A.** Simulations in network
 793 regime \mathbf{z}_1 (center) with simulations given connectivity perturbations in the negative direction of the sen-
 794 sitivity vector \mathbf{v}_1 (left) and positive direction (right). **B.** Same as A for network regime \mathbf{z}_2 figure.13). In
 795 regime 1, Anti accuracy diminishes in either direction of sensitivity away from the mode, while in
 796 regime 2, Anti accuracy tracks monotonic increases in Pro accuracy. These responses make intuitive
 797 sense, recognizing that \mathbf{v}_1 (unlike \mathbf{v}_2) points strongly in the direction of connectivity eigenvalue
 798 λ_{diag} , which is strongly anticorrelated with p_A (Fig. 14(SC3): **A.** Invariant eigenvectors of connectivity
 799 matrix W . **B.** Eigenvalues of connectivities of EPI distribution versus task accuraciesfigure.14, 15(SC4):
 800 **A.** Pairplots of eigenvalues of connectivity matrices in EPI distribution colored by Pro task accuracy. Red
 801 and purple stars and arrows correspond to eigenvalues and sensitivity directions \mathbf{z}_1 , \mathbf{z}_2 , \mathbf{v}_1 , and \mathbf{v}_2 . **B.** Same
 802 colored by Anti task accuracyfigure.15, see Section 5.2.4Superior colliculussubsubsection.5.2.4).
 803 In agreement with experimental results from Duan et al., we found optogenetic inactivation during

the delay period consistently decreased performance in the Anti task, but had no effect on the Pro task (Fig. 3A. 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 posterior distribution of rapid task switching networks. Red and purple stars (\mathbf{z}_1 and \mathbf{z}_2) indicate different connectivity regimes with different sensitivity vectors \mathbf{v}_1 and \mathbf{v}_2 . (Middle-left) Posterior predictive distribution of task accuracies. (Bottom-left) Task accuracy along dimensions of sensitivity in each connectivity regime. D. Means (solid) and standard deviations (shaded) of each population across random simulated trials. Top plots show Pro (top) and Anti (bottom) responses for connectivity \mathbf{z}_1 . Bottom rows show the same \mathbf{z}_2 . E. The EPI posterior predicts experimental results (left) showing no change in the Pro task, but larger error in the Anti task (right). F. Accuracy in the Anti task during delay period optogenetic inactivation $p_{A,opto}$ is strongly anticorrelated with accuracy in the Pro task. G. Accuracy with delay period inactivation along each connectivity regime's dimension of sensitivity (figure.3E)). This difference in resiliency across tasks to delay perturbation is a prediction made by the inferred EPI distribution, rather than an emergent property that was conditioned upon. Similarities across Pro and Anti trials in choice period responses following delay period inactivation (Fig. 17(SC6): A. Entropy throughout optimization. B. The emergent property statistic means and variances converge to their constraints at 20,000 iterations following the tenth augmented Lagrangian epoch (figure.17A) suggested that connectivity patterns inducing greater Pro task accuracy increase error in delay period inactivated Anti trials (Fig. 3A. 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 posterior distribution of rapid task switching networks. Red and purple stars (\mathbf{z}_1 and \mathbf{z}_2) indicate different connectivity regimes with different sensitivity vectors \mathbf{v}_1 and \mathbf{v}_2 . (Middle-left) Posterior predictive distribution of task accuracies. (Bottom-left) Task accuracy along dimensions of sensitivity in each connectivity regime. D. Means (solid) and standard deviations (shaded) of each population across random simulated trials. Top plots show Pro (top) and Anti (bottom) responses for connectivity \mathbf{z}_1 . Bottom rows show the same \mathbf{z}_2 . E. The EPI posterior predicts experimental results (left) showing no change in the Pro task, but larger error in the Anti task (right). F. Accuracy in the Anti task during delay period optogenetic inactivation $p_{A,opto}$ is strongly anticorrelated with accuracy in the Pro task. G. Accuracy with delay period inactivation along each connectivity regime's dimension of sensitivity (figure.3F)). The strong anticorrelation between p_P and $p_{A,opto}$ across posterior connectivities led to the following hypothesis about each connectivity regime: the sensitivity dimension of each regime decreases $p_{A,opto}$ irrespective of its effect on p_A , since both \mathbf{v}_1 and \mathbf{v}_2 increase p_P . Indeed, in regimes 1 and 2 where sensitivity dimensions elicit different responses in p_A , $p_{A,opto}$ decreases since the connectivity changes enhancing p_P exacerbate Anti trial error (Fig. 3A. Rapid task switching behavioral paradigm (see

837 text). **B.** Model of superior colliculus (SC). Neurons: LP - left pro, RP - right pro, LA - left anti, RA - right anti.
 838 Parameters: sW - self, hW - horizontal, vW -vertical, dW - diagonal weights. **C.** The EPI posterior distribution of
 839 rapid task switching networks. Red and purple stars (\mathbf{z}_1 and \mathbf{z}_2) indicate different connectivity regimes with different
 840 sensitivity vectors \mathbf{v}_1 and \mathbf{v}_2 . (Middle-left) Posterior predictive distribution of task accuracies. (Bottom-left) Task
 841 accuracy along dimensions of sensitivity in each connectivity regime. **D.** Means (solid) and standard deviations
 842 (shaded) of each population across random simulated trials. Top plots show Pro (top) and Anti (bottom) responses
 843 for connectivity \mathbf{z}_1 . Bottom rows show the same \mathbf{z}_2 . **E.** The EPI posterior predicts experimental results (left) showing
 844 no change in the Pro task, but larger error in the Anti task (right). **F.** Accuracy in the Anti task during delay period
 845 optogenetic inactivation $p_{A,\text{opto}}$ is strongly anticorrelated with accuracy in the Pro task. **G.** Accuracy with delay
 846 period inactivation along each connectivity regime's dimension of sensitivity figure.3F).

847 In summary, we used EPI to obtain the full distribution of connectivities that execute rapid task
 848 switching. This posterior revealed multiple regimes of rapid task switching, which we characterized
 849 using the probabilistic toolkit EPI seemlessly affords. EPI allowed us to conclude that since *all*
 850 parameters of this model producing rapid task switching make an experimentally verified prediction,
 851 the model is well-chosen in that regard. Finally, we used our knowledge about how \mathbf{z} governs $p_{A,\text{opto}}$
 852 to make accurate predictions about each identified regime of connectivity.

853 3.5 EPI scales well to high-dimensional parameter spaces

854 Here, we study the scalability of EPI in number of parameters $|\mathbf{z}|$ by inferring the connectivities
 855 of recurrent neural networks (RNNs, Fig. 4A. Recurrent neural network. **B.** EPI scales with z to
 856 high dimensions. Convergence definitions: EPI (blue) - satisfies all moment constraints, SNPE (orange)-
 857 produces at least $2/n_{\text{train}}$ parameter samples are in the bounds of emergent property (mean ± 0.5), and
 858 SMC-ABC (red) - 100 particles with $\epsilon < 0.5$ are produced. **C.** Posterior predictive distributions of EPI (blue),
 859 SNPE (orange), and SMC-ABC (red). Gray star indicates emergent property mean, and gray dashed lines
 860 indicate two standard deviations corresponding to the variance constraint. For $N \leq 6$ where SMC-ABC
 861 converges, samples are not diverse (path degeneracies). For $N \geq 25$, SNPE does not produce a posterior
 862 approximation yielding parameters with simulations near x_0 . **D.** Simulations of network parameters resulting
 863 from each method ($\tau = 100ms$). Each trace corresponds to simulation of one z . (Below) Ratio of obtained
 864 samples producing stable amplification figure.4A). We consider a rank-2 RNN with N neurons having
 865 connectivity

$$W = UV^\top + g\chi \quad (8)$$

866 and dynamics

$$\tau \dot{\mathbf{x}} = -\mathbf{x} + W\mathbf{x} \quad (9)$$

867 where $U = [\mathbf{u}_1 \ \mathbf{u}_2]$, $V = [\mathbf{v}_1 \ \mathbf{v}_2]$ and $\mathbf{u}_1, \mathbf{u}_2, \mathbf{v}_1, \mathbf{v}_2 \in [-1, 1]^N$. The random component has
 868 strength $g = 0.01$ and $\chi_{i,j} \sim \mathcal{N}(0, 1)$. We infer connectivity distributions $\mathbf{z} = [\mathbf{u}_1^\top, \mathbf{u}_2^\top, \mathbf{v}_1^\top, \mathbf{v}_2^\top]^\top$
 869 producing stable amplification. RNN's exhibiting stable amplification amplify responses to input
 870 along some dimensions, and are stable across all dimensions. Two conditions are both necessary
 871 and sufficient for RNNs to exhibit stable amplification [74]: $\text{real}(\lambda_1) < 1$ and $\lambda_1^s > 1$, where λ_1 is
 872 the eigenvalue of W with greatest real part and λ_1^s is the maximum eigenvalue of $W^s = \frac{W+W^\top}{2}$.

873 In our analysis, we seek to condition rank-2 networks of increasing size on a regime of stable ampli-
 874 fication. Networks with $\text{real}(\lambda_1) = 0.5 \pm 0.5$ and $\lambda_1^s = 1.5 \pm 0.5$ will yield moderate amplification.
 875 EPI can naturally condition on this emergent property

$$\begin{aligned} \mathcal{X} : \mathbb{E}_{\mathbf{z}, \mathbf{x}} \begin{bmatrix} \text{real}(\lambda_1) \\ \lambda_1^s \end{bmatrix} &= \begin{bmatrix} 0.5 \\ 1.5 \end{bmatrix} \\ \text{Var}_{\mathbf{z}, \mathbf{x}} \begin{bmatrix} \text{real}(\lambda_1) \\ \lambda_1^s \end{bmatrix} &= \begin{bmatrix} 0.25^2 \\ 0.25^2 \end{bmatrix}. \end{aligned} \quad (10)$$

876 For comparison, we infer rank-2 RNN connectivities with alternative approaches to likelihood free-
 877 inference. ABC methods define a tolerance ϵ from observed data x_0 for which we keep sampled
 878 parameters. To make this ABC approach as similar as possible to the EPI program defined by
 879 Equation 5EPI scales well to high-dimensional parameter spacesequation.3.5, we chose $\epsilon = 0.5$, an
 880 l_2 distance metric, and

$$x_0 = \begin{bmatrix} \text{real}(\lambda_1) \\ \lambda_1^s \end{bmatrix} = \begin{bmatrix} 0.5 \\ 1.5 \end{bmatrix} \quad (11)$$

881 located at the mean of our desired emergent property. We use sequential Monte Carlo ABC (SMC-
 882 ABC), to compare efficiency, since it is considered the most efficient ABC approach. SNPE [75] is
 883 another deep likelihood-free inference method that emerged along with this work. In contrast to
 884 EPI, SNPE cannot condition on the variance of the posterior predictive distribution. Also, there
 885 is no tolerance parameter for SNPE like ϵ in ABC, so the comparative SNPE approach simply
 886 conditions on observation x_0 .

887 As we scale the number of neurons N in the RNN, and thus the dimensionality of the parameter
 888 space $\mathbf{z} \in [-1, 1]^{4N}$, we see that EPI has superior scaling properties (Fig. 4A. Recurrent neural
 889 network. **B.** EPI scales with z to high dimensions. Convergence definitions: EPI (blue) - satisfies all

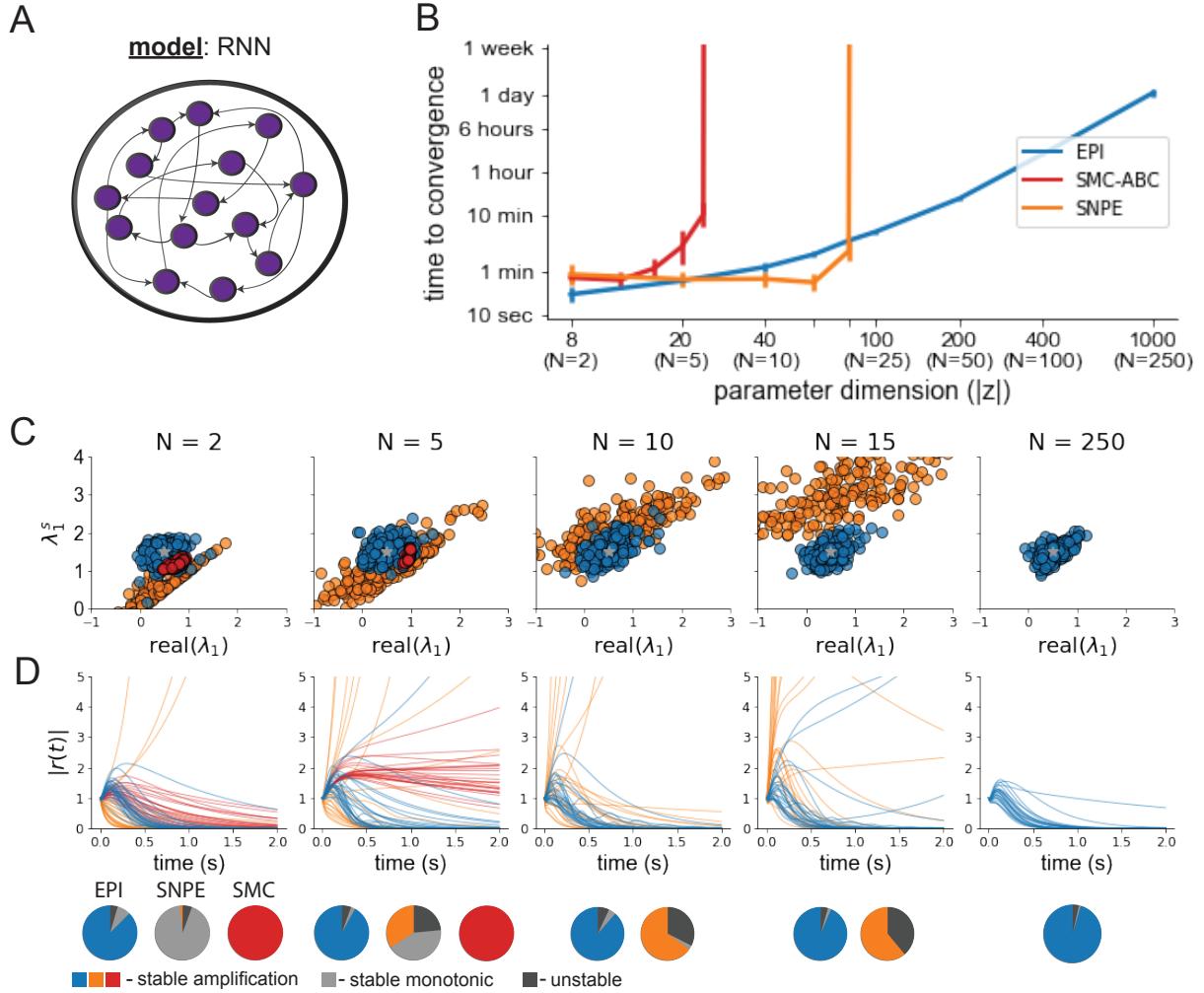


Figure 4: **A.** Recurrent neural network. **B.** EPI scales with z to high dimensions. Convergence definitions: EPI (blue) - satisfies all moment constraints, SNPE (orange)- produces at least $2/n_{\text{train}}$ parameter samples are in the bounds of emergent property (mean ± 0.5), and SMC-ABC (red) - 100 particles with $\epsilon < 0.5$ are produced. **C.** Posterior predictive distributions of EPI (blue), SNPE (orange), and SMC-ABC (red). Gray star indicates emergent property mean, and gray dashed lines indicate two standard deviations corresponding to the variance constraint. For $N \leq 6$ where SMC-ABC converges, samples are not diverse (path degeneracies). For $N \geq 25$, SNPE does not produce a posterior approximation yielding parameters with simulations near x_0 . **D.** Simulations of network parameters resulting from each method ($\tau = 100ms$). Each trace corresponds to simulation of one z . (Below) Ratio of obtained samples producing stable amplification.

890 moment constraints, SNPE (orange)- produces at least $2/n_{\text{train}}$ parameter samples are in the bounds of
 891 emergent property (mean \pm 0.5), and SMC-ABC (red) - 100 particles with $\epsilon < 0.5$ are produced. **C.**
 892 Posterior predictive distributions of EPI (blue), SNPE (orange), and SMC-ABC (red). Gray star indicates
 893 emergent property mean, and gray dashed lines indicate two standard deviations corresponding to the
 894 variance constraint. For $N \leq 6$ where SMC-ABC converges, samples are not diverse (path degeneracies).
 895 For $N \geq 25$, SNPE does not produce a posterior approximation yielding parameters with simulations near
 896 x_0 . **D.** Simulations of network parameters resulting from each method ($\tau = 100ms$). Each trace corresponds
 897 to simulation of one z . (Below) Ratio of obtained samples producing stable amplification figure.4B). SMC-
 898 ABC and SNPE become intractable around 25 and 90 dimensions respectively, while EPI can infer
 899 1000-dimensional distributions in about 1 day. No matter the number of neurons, EPI always
 900 produces connectivity distributions with mean and variance of $\text{real}(\lambda_1)$ and λ_1^s of \mathcal{X} (Fig. 4A.
 901 Recurrent neural network. **B.** EPI scales with z to high dimensions. Convergence definitions: EPI (blue)
 902 - satisfies all moment constraints, SNPE (orange)- produces at least $2/n_{\text{train}}$ parameter samples are in the
 903 bounds of emergent property (mean \pm 0.5), and SMC-ABC (red) - 100 particles with $\epsilon < 0.5$ are produced.
 904 **C.** Posterior predictive distributions of EPI (blue), SNPE (orange), and SMC-ABC (red). Gray star indicates
 905 emergent property mean, and gray dashed lines indicate two standard deviations corresponding to the
 906 variance constraint. For $N \leq 6$ where SMC-ABC converges, samples are not diverse (path degeneracies).
 907 For $N \geq 25$, SNPE does not produce a posterior approximation yielding parameters with simulations near
 908 x_0 . **D.** Simulations of network parameters resulting from each method ($\tau = 100ms$). Each trace corresponds
 909 to simulation of one z . (Below) Ratio of obtained samples producing stable amplification figure.4C, blue),
 910 and high variation in response profiles 4A. Recurrent neural network. **B.** EPI scales with z to high
 911 dimensions. Convergence definitions: EPI (blue) - satisfies all moment constraints, SNPE (orange)- produces
 912 at least $2/n_{\text{train}}$ parameter samples are in the bounds of emergent property (mean \pm 0.5), and SMC-
 913 ABC (red) - 100 particles with $\epsilon < 0.5$ are produced. **C.** Posterior predictive distributions of EPI (blue),
 914 SNPE (orange), and SMC-ABC (red). Gray star indicates emergent property mean, and gray dashed lines
 915 indicate two standard deviations corresponding to the variance constraint. For $N \leq 6$ where SMC-ABC
 916 converges, samples are not diverse (path degeneracies). For $N \geq 25$, SNPE does not produce a posterior
 917 approximation yielding parameters with simulations near x_0 . **D.** Simulations of network parameters resulting
 918 from each method ($\tau = 100ms$). Each trace corresponds to simulation of one z . (Below) Ratio of obtained
 919 samples producing stable amplification figure.4D, blue). For the dimensionalities in which SMC-ABC is
 920 tractable, the inferred parameters always exhibit stable amplification, are less varied 4A. Recurrent
 921 neural network. **B.** EPI scales with z to high dimensions. Convergence definitions: EPI (blue) - satisfies
 922 all moment constraints, SNPE (orange)- produces at least $2/n_{\text{train}}$ parameter samples are in the bounds

923 of emergent property (mean \pm 0.5), and SMC-ABC (red) - 100 particles with $\epsilon < 0.5$ are produced. **C.**
 924 Posterior predictive distributions of EPI (blue), SNPE (orange), and SMC-ABC (red). Gray star indicates
 925 emergent property mean, and gray dashed lines indicate two standard deviations corresponding to the
 926 variance constraint. For $N \leq 6$ where SMC-ABC converges, samples are not diverse (path degeneracies).
 927 For $N \geq 25$, SNPE does not produce a posterior approximation yielding parameters with simulations near
 928 x_0 . **D.** Simulations of network parameters resulting from each method ($\tau = 100ms$). Each trace corresponds
 929 to simulation of one z . (Below) Ratio of obtained samples producing stable amplificationfigure.4C, red) and
 930 largely produce similar responses **4A.** Recurrent neural network. **B.** EPI scales with z to high dimensions.
 931 Convergence definitions: EPI (blue) - satisfies all moment constraints, SNPE (orange)- produces at least
 932 $2/n_{\text{train}}$ parameter samples are in the bounds of emergent property (mean \pm 0.5), and SMC-ABC (red) - 100
 933 particles with $\epsilon < 0.5$ are produced. **C.** Posterior predictive distributions of EPI (blue), SNPE (orange), and
 934 SMC-ABC (red). Gray star indicates emergent property mean, and gray dashed lines indicate two standard
 935 deviations corresponding to the variance constraint. For $N \leq 6$ where SMC-ABC converges, samples are
 936 not diverse (path degeneracies). For $N \geq 25$, SNPE does not produce a posterior approximation yielding
 937 parameters with simulations near x_0 . **D.** Simulations of network parameters resulting from each method
 938 ($\tau = 100ms$). Each trace corresponds to simulation of one z . (Below) Ratio of obtained samples producing
 939 stable amplificationfigure.4D, red). When using SNPE the inferred parameters are widely varied **4A.**
 940 Recurrent neural network. **B.** EPI scales with z to high dimensions. Convergence definitions: EPI (blue)
 941 - satisfies all moment constraints, SNPE (orange)- produces at least $2/n_{\text{train}}$ parameter samples are in the
 942 bounds of emergent property (mean \pm 0.5), and SMC-ABC (red) - 100 particles with $\epsilon < 0.5$ are produced.
 943 **C.** Posterior predictive distributions of EPI (blue), SNPE (orange), and SMC-ABC (red). Gray star indicates
 944 emergent property mean, and gray dashed lines indicate two standard deviations corresponding to the
 945 variance constraint. For $N \leq 6$ where SMC-ABC converges, samples are not diverse (path degeneracies).
 946 For $N \geq 25$, SNPE does not produce a posterior approximation yielding parameters with simulations near
 947 x_0 . **D.** Simulations of network parameters resulting from each method ($\tau = 100ms$). Each trace corresponds
 948 to simulation of one z . (Below) Ratio of obtained samples producing stable amplificationfigure.4C, orange),
 949 but often produce non-amplified or unstable responses **4A.** Recurrent neural network. **B.** EPI scales
 950 with z to high dimensions. Convergence definitions: EPI (blue) - satisfies all moment constraints, SNPE
 951 (orange)- produces at least $2/n_{\text{train}}$ parameter samples are in the bounds of emergent property (mean \pm
 952 0.5), and SMC-ABC (red) - 100 particles with $\epsilon < 0.5$ are produced. **C.** Posterior predictive distributions
 953 of EPI (blue), SNPE (orange), and SMC-ABC (red). Gray star indicates emergent property mean, and
 954 gray dashed lines indicate two standard deviations corresponding to the variance constraint. For $N \leq 6$
 955 where SMC-ABC converges, samples are not diverse (path degeneracies). For $N \geq 25$, SNPE does not

956 produce a posterior approximation yielding parameters with simulations near x_0 . **D.** Simulations of network
957 parameters resulting from each method ($\tau = 100ms$). Each trace corresponds to simulation of one z .
958 (Below) Ratio of obtained samples producing stable amplificationfigure.4D, orange). In conclusion, we
959 found that deep likelihood-free inference techniques are capable of scaling to higher dimensional
960 inference than SMC-ABC. However, only EPI can scale to high dimensions while reproducing the
961 emergent property.

962 4 Discussion

963 In neuroscience, machine learning has primarily been used to reveal structure in neural datasets
964 [30]. Such careful inference procedures are developed for these statistical models allowing precise,
965 quantitative reasoning, which clarifies the way data informs beliefs about the model parameters.
966 However, these statistical models lack resemblance to the underlying biology, making it unclear
967 how to go from the structure revealed by these methods, to the neural mechanisms giving rise
968 to it. In contrast, theoretical neuroscience has focused on careful mechanistic modeling and the
969 production of emergent properties of computation. The careful steps of *i.*) model design and
970 *ii.*) emergent property definition, are followed by *iii.*) practical inference methods resulting in an
971 opaque characterization of the way model parameters govern computation. In this work, we improve
972 upon parameter inference techniques in theoretical neuroscience with emergent property inference,
973 harnessing deep learning towards careful inference in careful models of neural computation (see
974 Section 5.1.1Related approachessubsubsection.5.1.1).

975 Specifically, approximate Bayesian computation [42, 43, 31] has been the standard approach to
976 parameter inference in neural circuit models lacking tractable likelihoods. ABC methods do not
977 confer probabilities on accepted parameters, require an acceptance threshold chosen to trade-off
978 inference quality with tractability, do not scale efficiently to high-dimensional parameter spaces, and
979 require independent techniques to analyze sensitivity for local parameter choices [76]. In contrast,
980 EPI allows probability evaluations at any point in parameter space, conditions posteriors on the
981 natural quantification of emergent properties, scales to high dimensional parameter spaces, and
982 naturally admits sensitivity quantification via fast evaluations of the posterior Hessian.

983 Technically, EPI is a maximum entropy method, which learns parameter distributions that are
984 as random as possible given that they produce the emergent property. Conceptually, maximally
985 random distributions given some constraints are useful for understanding parametric sensitivity.

986 This is well understood in Bayesian inference, where maximum entropy is the chosen normative
987 principle. This is emphasized by an innovative formalism unifying top-down maximum entropy
988 normative models with bottom-up statistical models [77]. Indeed, EPI is an adaptive variational
989 inference program, and may be considered to have a Bayesian uniform prior (see Section 5.1.6EPI
990 as variational inferencesubsubsection.5.1.6).

991 Biologically realistic models of neural circuits often prove formidable to analyze for two main
992 reasons. A primary challenge is that the number of parameters scales dramatically with the number
993 of neurons, limiting analysis of its parameter space. We see in Section 3.5EPI scales well to
994 high-dimensional parameter spacessubsection.3.5 that EPI scales seemlessly to high dimensional
995 parameter spaces of RNN connectivities, while maintaining the production of the specified emergent
996 property. EPI strongly outperforms the standard likelihood-free inference technique (SMC-ABC
997 [31]), and a recently developed deep likelihood-free inference technique (SNPE [75]), most likely
998 because of it's ability to leverage the gradient information of the emergent property statistics and
999 to adapt it's paramter sampling distribution at every step of gradient descent.

1000 A secondary challenge is that the structure of the parametric regimes governing emergent properties
1001 is intricate. For example, even in low dimensional circuits, models can support more than one steady
1002 state [78] and non-trivial dynamics on strange attractors [79]. With EPI, we use deep probabillity
1003 distributions to capture the complex nonlinear parameter distributions governing model behavior.

1004 In Section 3.3EPI reveals how neuron-type specific noise governs variability in the stochastic sta-
1005 bilized supralinear networkssubsection.3.3, we used EPI to reveal a curved parametric manifolds
1006 governing circuit variability in the stochastic stabilized supralinear network, and used hypothesis
1007 testing techniques to validate our findings. In Section 3.4EPI identifies multiple regimes of rapid
1008 task switchingsubsection.3.4, we identified two regimes of SC connectivity resulting in rapid task
1009 switching, and found that the full distribution of rapid task switching networks reproduced an
1010 experimental result.

1011 EPI leverages deep learning technology for neuroscientific inquiry in a categorically different way
1012 than approaches focused on training neural networks to execute behavioral tasks [80]. These works
1013 focus on examining optimized deep neural networks while considering the objective function, learn-
1014 ing rule, and architecture used. This endeavor efficiently obtains sets of parameters that can be
1015 reasoned about with respect to such considerations, but lacks the careful probabilistic treatment of
1016 parameter inference in EPI. All of these approaches can be used complementarily to enhance the
1017 practice of theoretical neuroscience.

1018 **Acknowledgements:**

1019 This work was funded by NSF Graduate Research Fellowship, DGE-1644869, McKnight Endow-
1020 ment Fund, NIH NINDS 5R01NS100066, Simons Foundation 542963, NSF NeuroNex Award, DBI-
1021 1707398, The Gatsby Charitable Foundation, Simons Collaboration on the Global Brain Postdoc-
1022 toral Fellowship, Chinese Postdoctoral Science Foundation, and International Exchange Program
1023 Fellowship. Helpful conversations were had with Francesca Mastrogiuseppe, Srdjan Ostojic, James
1024 Fitzgerald, Stephen Baccus, Dhruva Raman, Liam Paninski, and Larry Abbott.

1025 **Data availability statement:**

1026 The datasets generated during and/or analyzed during the current study are available from the
1027 corresponding author upon reasonable request.

1028 **Code availability statement:**

1029 All software written for the current study is available at <https://github.com/cunningham-lab/epi>.

1030 **References**

- 1031 [1] Nancy Kopell and G Bard Ermentrout. Coupled oscillators and the design of central pattern
1032 generators. *Mathematical biosciences*, 90(1-2):87–109, 1988.
- 1033 [2] Eve Marder. From biophysics to models of network function. *Annual review of neuroscience*,
1034 21(1):25–45, 1998.
- 1035 [3] Larry F Abbott. Theoretical neuroscience rising. *Neuron*, 60(3):489–495, 2008.
- 1036 [4] Xiao-Jing Wang. Neurophysiological and computational principles of cortical rhythms in
1037 cognition. *Physiological reviews*, 90(3):1195–1268, 2010.
- 1038 [5] Ryan N Gutenkunst, Joshua J Waterfall, Fergal P Casey, Kevin S Brown, Christopher R
1039 Myers, and James P Sethna. Universally sloppy parameter sensitivities in systems biology
1040 models. *PLoS Comput Biol*, 3(10):e189, 2007.
- 1041 [6] Timothy O’Leary, Alex H Williams, Alessio Franci, and Eve Marder. Cell types, network
1042 homeostasis, and pathological compensation from a biologically plausible ion channel expres-
1043 sion model. *Neuron*, 82(4):809–821, 2014.
- 1044 [7] John J Hopfield. Neural networks and physical systems with emergent collective computa-
1045 tional abilities. *Proceedings of the national academy of sciences*, 79(8):2554–2558, 1982.

- 1046 [8] Haim Sompolinsky, Andrea Crisanti, and Hans-Jurgen Sommers. Chaos in random neural
1047 networks. *Physical review letters*, 61(3):259, 1988.
- 1048 [9] Misha V Tsodyks, William E Skaggs, Terrence J Sejnowski, and Bruce L McNaughton. Para-
1049 doxical effects of external modulation of inhibitory interneurons. *Journal of neuroscience*,
1050 17(11):4382–4388, 1997.
- 1051 [10] Kong-Fatt Wong and Xiao-Jing Wang. A recurrent network mechanism of time integration
1052 in perceptual decisions. *Journal of Neuroscience*, 26(4):1314–1328, 2006.
- 1053 [11] WR Foster, LH Ungar, and JS Schwaber. Significance of conductances in hodgkin-huxley
1054 models. *Journal of neurophysiology*, 70(6):2502–2518, 1993.
- 1055 [12] Astrid A Prinz, Dirk Bucher, and Eve Marder. Similar network activity from disparate circuit
1056 parameters. *Nature neuroscience*, 7(12):1345–1352, 2004.
- 1057 [13] Pablo Achard and Erik De Schutter. Complex parameter landscape for a complex neuron
1058 model. *PLoS computational biology*, 2(7):e94, 2006.
- 1059 [14] Leandro M Alonso and Eve Marder. Visualization of currents in neural models with similar
1060 behavior and different conductance densities. *Elife*, 8:e42722, 2019.
- 1061 [15] Robert E Kass and Valérie Ventura. A spike-train probability model. *Neural computation*,
1062 13(8):1713–1720, 2001.
- 1063 [16] Emery N Brown, Loren M Frank, Dengda Tang, Michael C Quirk, and Matthew A Wilson.
1064 A statistical paradigm for neural spike train decoding applied to position prediction from
1065 ensemble firing patterns of rat hippocampal place cells. *Journal of Neuroscience*, 18(18):7411–
1066 7425, 1998.
- 1067 [17] Liam Paninski. Maximum likelihood estimation of cascade point-process neural encoding
1068 models. *Network: Computation in Neural Systems*, 15(4):243–262, 2004.
- 1069 [18] Wilson Truccolo, Uri T Eden, Matthew R Fellows, John P Donoghue, and Emery N Brown.
1070 A point process framework for relating neural spiking activity to spiking history, neural
1071 ensemble, and extrinsic covariate effects. *Journal of neurophysiology*, 93(2):1074–1089, 2005.
- 1072 [19] Elad Schneidman, Michael J Berry, Ronen Segev, and William Bialek. Weak pairwise corre-
1073 lations imply strongly correlated network states in a neural population. *Nature*, 440(7087):1007–
1074 1012, 2006.

- 1075 [20] Shaul Druckmann, Yoav Banitt, Albert A Gidon, Felix Schürmann, Henry Markram, and Idan
1076 Segev. A novel multiple objective optimization framework for constraining conductance-based
1077 neuron models by experimental data. *Frontiers in neuroscience*, 1:1, 2007.
- 1078 [21] Richard Turner and Maneesh Sahani. A maximum-likelihood interpretation for slow feature
1079 analysis. *Neural computation*, 19(4):1022–1038, 2007.
- 1080 [22] M Yu Byron, John P Cunningham, Gopal Santhanam, Stephen I Ryu, Krishna V Shenoy, and
1081 Maneesh Sahani. Gaussian-process factor analysis for low-dimensional single-trial analysis of
1082 neural population activity. In *Advances in neural information processing systems*, pages
1083 1881–1888, 2009.
- 1084 [23] Jakob H Macke, Lars Buesing, John P Cunningham, Byron M Yu, Krishna V Shenoy, and
1085 Maneesh Sahani. Empirical models of spiking in neural populations. *Advances in neural*
1086 *information processing systems*, 24:1350–1358, 2011.
- 1087 [24] Il Memming Park and Jonathan W Pillow. Bayesian spike-triggered covariance analysis. In
1088 *Advances in neural information processing systems*, pages 1692–1700, 2011.
- 1089 [25] Einat Granot-Atedgi, Gašper Tkačik, Ronen Segev, and Elad Schneidman. Stimulus-
1090 dependent maximum entropy models of neural population codes. *PLoS Comput Biol*,
1091 9(3):e1002922, 2013.
- 1092 [26] Kenneth W Latimer, Jacob L Yates, Miriam LR Meister, Alexander C Huk, and Jonathan W
1093 Pillow. Single-trial spike trains in parietal cortex reveal discrete steps during decision-making.
1094 *Science*, 349(6244):184–187, 2015.
- 1095 [27] Kaushik J Lakshminarasimhan, Marina Petsalis, Hyeshin Park, Gregory C DeAngelis, Xaq
1096 Pitkow, and Dora E Angelaki. A dynamic bayesian observer model reveals origins of bias in
1097 visual path integration. *Neuron*, 99(1):194–206, 2018.
- 1098 [28] Lea Duncker, Gergo Bohner, Julien Boussard, and Maneesh Sahani. Learning interpretable
1099 continuous-time models of latent stochastic dynamical systems. *Proceedings of the 36th In-*
1100 *ternational Conference on Machine Learning*, 2019.
- 1101 [29] Josef Ladenbauer, Sam McKenzie, Daniel Fine English, Olivier Hagens, and Srdjan Ostojic.
1102 Inferring and validating mechanistic models of neural microcircuits based on spike-train data.
1103 *Nature Communications*, 10(4933), 2019.

- 1104 [30] Liam Paninski and John P Cunningham. Neural data science: accelerating the experiment-
1105 analysis-theory cycle in large-scale neuroscience. *Current opinion in neurobiology*, 50:232–241,
1106 2018.
- 1107 [31] Scott A Sisson, Yanan Fan, and Mark M Tanaka. Sequential monte carlo without likelihoods.
1108 *Proceedings of the National Academy of Sciences*, 104(6):1760–1765, 2007.
- 1109 [32] Juliane Liepe, Paul Kirk, Sarah Filippi, Tina Toni, Chris P Barnes, and Michael PH Stumpf.
1110 A framework for parameter estimation and model selection from experimental data in systems
1111 biology using approximate bayesian computation. *Nature protocols*, 9(2):439–456, 2014.
- 1112 [33] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *International Con-*
1113 *ference on Learning Representations*, 2014.
- 1114 [34] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropaga-
1115 tion and variational inference in deep latent gaussian models. *International Conference on*
1116 *Machine Learning*, 2014.
- 1117 [35] Yuanjun Gao, Evan W Archer, Liam Paninski, and John P Cunningham. Linear dynamical
1118 neural population models through nonlinear embeddings. In *Advances in neural information*
1119 *processing systems*, pages 163–171, 2016.
- 1120 [36] Yuan Zhao and Il Memming Park. Recursive variational bayesian dual estimation for non-
1121 linear dynamics and non-gaussian observations. *stat*, 1050:27, 2017.
- 1122 [37] Gabriel Barello, Adam Charles, and Jonathan Pillow. Sparse-coding variational auto-
1123 encoders. *bioRxiv*, page 399246, 2018.
- 1124 [38] Chethan Pandarinath, Daniel J O’Shea, Jasmine Collins, Rafal Jozefowicz, Sergey D Stavisky,
1125 Jonathan C Kao, Eric M Trautmann, Matthew T Kaufman, Stephen I Ryu, Leigh R
1126 Hochberg, et al. Inferring single-trial neural population dynamics using sequential auto-
1127 encoders. *Nature methods*, page 1, 2018.
- 1128 [39] Alexander B Wiltschko, Matthew J Johnson, Giuliano Iurilli, Ralph E Peterson, Jesse M
1129 Katon, Stan L Pashkovski, Victoria E Abraira, Ryan P Adams, and Sandeep Robert Datta.
1130 Mapping sub-second structure in mouse behavior. *Neuron*, 88(6):1121–1135, 2015.

- 1131 [40] Matthew J Johnson, David K Duvenaud, Alex Wiltschko, Ryan P Adams, and Sandeep R
1132 Datta. Composing graphical models with neural networks for structured representations and
1133 fast inference. In *Advances in neural information processing systems*, pages 2946–2954, 2016.
- 1134 [41] Eleanor Batty, Matthew Whiteway, Shreya Saxena, Dan Biderman, Taiga Abe, Simon Musall,
1135 Winthrop Gillis, Jeffrey Markowitz, Anne Churchland, John Cunningham, et al. Behavenet:
1136 nonlinear embedding and bayesian neural decoding of behavioral videos. *Advances in Neural*
1137 *Information Processing Systems*, 2019.
- 1138 [42] Mark A Beaumont, Wenyang Zhang, and David J Balding. Approximate bayesian computa-
1139 tion in population genetics. *Genetics*, 162(4):2025–2035, 2002.
- 1140 [43] Paul Marjoram, John Molitor, Vincent Plagnol, and Simon Tavaré. Markov chain monte carlo
1141 without likelihoods. *Proceedings of the National Academy of Sciences*, 100(26):15324–15328,
1142 2003.
- 1143 [44] W Keith Hastings. Monte carlo sampling methods using markov chains and their applications.
1144 1970.
- 1145 [45] Nicholas Metropolis, Arianna W Rosenbluth, Marshall N Rosenbluth, Augusta H Teller, and
1146 Edward Teller. Equation of state calculations by fast computing machines. *The journal of*
1147 *chemical physics*, 21(6):1087–1092, 1953.
- 1148 [46] Lawrence Saul and Michael Jordan. A mean field learning algorithm for unsupervised neural
1149 networks. In *Learning in graphical models*, pages 541–554. Springer, 1998.
- 1150 [47] Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows.
1151 *International Conference on Machine Learning*, 2015.
- 1152 [48] Mark K Transtrum, Benjamin B Machta, Kevin S Brown, Bryan C Daniels, Christopher R
1153 Myers, and James P Sethna. Perspective: Sloppiness and emergent theories in physics,
1154 biology, and beyond. *The Journal of chemical physics*, 143(1):07B201_1, 2015.
- 1155 [49] Dustin Tran, Rajesh Ranganath, and David Blei. Hierarchical implicit models and likelihood-
1156 free variational inference. In *Advances in Neural Information Processing Systems*, pages
1157 5523–5533, 2017.
- 1158 [50] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp.
1159 *Proceedings of the 5th International Conference on Learning Representations*, 2017.

- 1160 [51] George Papamakarios, Theo Pavlakou, and Iain Murray. Masked autoregressive flow for
1161 density estimation. In *Advances in Neural Information Processing Systems*, pages 2338–2347,
1162 2017.
- 1163 [52] Gabriel Loaiza-Ganem, Yuanjun Gao, and John P Cunningham. Maximum entropy flow
1164 networks. *International Conference on Learning Representations*, 2017.
- 1165 [53] Mark S Goldman, Jorge Golowasch, Eve Marder, and LF Abbott. Global structure, ro-
1166 bustness, and modulation of neuronal models. *Journal of Neuroscience*, 21(14):5229–5238,
1167 2001.
- 1168 [54] Ashok Litwin-Kumar, Robert Rosenbaum, and Brent Doiron. Inhibitory stabilization and
1169 visual coding in cortical circuits with multiple interneuron subtypes. *Journal of neurophysi-
1170 ology*, 115(3):1399–1409, 2016.
- 1171 [55] Chunyu A Duan, Marino Pagan, Alex T Piet, Charles D Kopec, Athena Akrami, Alexander J
1172 Riordan, Jeffrey C Erlich, and Carlos D Brody. Collicular circuits for flexible sensorimotor
1173 routing. *bioRxiv*, page 245613, 2018.
- 1174 [56] Eve Marder and Vatsala Thirumalai. Cellular, synaptic and network effects of neuromodula-
1175 tion. *Neural Networks*, 15(4-6):479–493, 2002.
- 1176 [57] Gabrielle J Gutierrez, Timothy O’Leary, and Eve Marder. Multiple mechanisms switch an
1177 electrically coupled, synaptically inhibited neuron between competing rhythmic oscillators.
1178 *Neuron*, 77(5):845–858, 2013.
- 1179 [58] Edwin T Jaynes. Information theory and statistical mechanics. *Physical review*, 106(4):620,
1180 1957.
- 1181 [59] Gamaleldin F Elsayed and John P Cunningham. Structure in neural population recordings:
1182 an expected byproduct of simpler phenomena? *Nature neuroscience*, 20(9):1310, 2017.
- 1183 [60] Cristina Savin and Gašper Tkačik. Maximum entropy models as a tool for building precise
1184 neural controls. *Current opinion in neurobiology*, 46:120–126, 2017.
- 1185 [61] Mark S Goldman. Memory without feedback in a neural network. *Neuron*, 61(4):621–634,
1186 2009.
- 1187 [62] Brendan K Murphy and Kenneth D Miller. Balanced amplification: a new mechanism of
1188 selective amplification of neural activity patterns. *Neuron*, 61(4):635–648, 2009.

- 1189 [63] Hirofumi Ozeki, Ian M Finn, Evan S Schaffer, Kenneth D Miller, and David Ferster. Inhibitory
1190 stabilization of the cortical network underlies visual surround suppression. *Neuron*, 62(4):578–
1191 592, 2009.
- 1192 [64] Daniel B Rubin, Stephen D Van Hooser, and Kenneth D Miller. The stabilized supralinear-
1193 ear network: a unifying circuit motif underlying multi-input integration in sensory cortex.
1194 *Neuron*, 85(2):402–417, 2015.
- 1195 [65] Henry Markram, Maria Toledo-Rodriguez, Yun Wang, Anirudh Gupta, Gilad Silberberg, and
1196 Caizhi Wu. Interneurons of the neocortical inhibitory system. *Nature reviews neuroscience*,
1197 5(10):793, 2004.
- 1198 [66] Bernardo Rudy, Gordon Fishell, SooHyun Lee, and Jens Hjerling-Leffler. Three groups of
1199 interneurons account for nearly 100% of neocortical gabaergic neurons. *Developmental neu-
1200 robiology*, 71(1):45–61, 2011.
- 1201 [67] Robin Tremblay, Soohyun Lee, and Bernardo Rudy. GABAergic Interneurons in the Neocor-
1202 tex: From Cellular Properties to Circuits. *Neuron*, 91(2):260–292, 2016.
- 1203 [68] Carsten K Pfeffer, Mingshan Xue, Miao He, Z Josh Huang, and Massimo Scanziani. Inhi-
1204 bition of inhibition in visual cortex: the logic of connections between molecularly distinct
1205 interneurons. *Nature Neuroscience*, 16(8):1068, 2013.
- 1206 [69] Luis Carlos Garcia Del Molino, Guangyu Robert Yang, Jorge F. Mejias, and Xiao Jing
1207 Wang. Paradoxical response reversal of top- down modulation in cortical circuits with three
1208 interneuron types. *Elife*, 6:1–15, 2017.
- 1209 [70] Guang Chen, Carl Van Vreeswijk, David Hansel, and David Hansel. Mechanisms underlying
1210 the response of mouse cortical networks to optogenetic manipulation. 2019.
- 1211 [71] Guillaume Hennequin, Yashar Ahmadian, Daniel B Rubin, Máté Lengyel, and Kenneth D
1212 Miller. The dynamical regime of sensory cortex: stable dynamics around a single stimulus-
1213 tuned attractor account for patterns of noise variability. *Neuron*, 98(4):846–860, 2018.
- 1214 [72] Agostina Palmigiano, Francesco Fumarola, Daniel P Mossing, Nataliya Kraynyukova, Hillel
1215 Adesnik, and Kenneth Miller. Structure and variability of optogenetic responses identify the
1216 operating regime of cortex. *bioRxiv*, 2020.

- 1217 [73] Chunyu A Duan, Jeffrey C Erlich, and Carlos D Brody. Requirement of prefrontal and
1218 midbrain regions for rapid executive control of behavior in the rat. *Neuron*, 86(6):1491–1503,
1219 2015.
- 1220 [74] Giulio Bondanelli and Srdjan Ostojic. Coding with transient trajectories in recurrent neural
1221 networks. *PLoS computational biology*, 16(2):e1007655, 2020.
- 1222 [75] Pedro J Gonçalves, Jan-Matthis Lueckmann, Michael Deistler, Marcel Nonnenmacher, Kaan
1223 Öcal, Giacomo Bassetto, Chaitanya Chintaluri, William F Podlaski, Sara A Haddad, Tim P
1224 Vogels, et al. Training deep neural density estimators to identify mechanistic models of neural
1225 dynamics. *bioRxiv*, page 838383, 2019.
- 1226 [76] Scott A Sisson, Yanan Fan, and Mark Beaumont. *Handbook of approximate Bayesian com-*
1227 *putation*. CRC Press, 2018.
- 1228 [77] Wiktor Młynarski, Michal Hledík, Thomas R Sokolowski, and Gašper Tkačik. Statistical
1229 analysis and optimality of neural systems. *bioRxiv*, page 848374, 2020.
- 1230 [78] Nataliya Kraynyukova and Tatjana Tchumatchenko. Stabilized supralinear network can give
1231 rise to bistable, oscillatory, and persistent activity. *Proceedings of the National Academy of*
1232 *Sciences*, 115(13):3464–3469, 2018.
- 1233 [79] Katherine Morrison, Anda Degeratu, Vladimir Itskov, and Carina Curto. Diversity of emer-
1234 gent dynamics in competitive threshold-linear networks: a preliminary report. *arXiv preprint*
1235 *arXiv:1605.04463*, 2016.
- 1236 [80] Blake A Richards and et al. A deep learning framework for neuroscience. *Nature Neuroscience*,
1237 2019.
- 1238 [81] Mark Girolami and Ben Calderhead. Riemann manifold langevin and hamiltonian monte
1239 carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*,
1240 73(2):123–214, 2011.
- 1241 [82] Andrew Golightly and Darren J Wilkinson. Bayesian parameter inference for stochastic bio-
1242 chemical network models using particle markov chain monte carlo. *Interface focus*, 1(6):807–
1243 820, 2011.
- 1244 [83] Kyle Cranmer, Johann Brehmer, and Gilles Louppe. The frontier of simulation-based infer-
1245 ence. *Proceedings of the National Academy of Sciences*, 2020.

- 1246 [84] Sean R Bittner, Agostina Palmigiano, Kenneth D Miller, and John P Cunningham. Degener-
1247 ate solution networks for theoretical neuroscience. *Computational and Systems Neuroscience*
1248 *Meeting (COSYNE), Lisbon, Portugal*, 2019.
- 1249 [85] Sean R Bittner, Alex T Piet, Chunyu A Duan, Agostina Palmigiano, Kenneth D Miller,
1250 Carlos D Brody, and John P Cunningham. Examining models in theoretical neuroscience
1251 with degenerate solution networks. *Bernstein Conference 2019, Berlin, Germany*, 2019.
- 1252 [86] Marcel Nonnenmacher, Pedro J Goncalves, Giacomo Bassetto, Jan-Matthis Lueckmann, and
1253 Jakob H Macke. Robust statistical inference for simulation-based models in neuroscience. In
1254 *Bernstein Conference 2018, Berlin, Germany*, 2018.
- 1255 [87] Deistler Michael, , Pedro J Goncalves, Kaan Oecal, and Jakob H Macke. Statistical infer-
1256 ence for analyzing sloppiness in neuroscience models. In *Bernstein Conference 2019, Berlin,*
1257 *Germany*, 2019.
- 1258 [88] Jan-Matthis Lueckmann, Pedro J Goncalves, Giacomo Bassetto, Kaan Öcal, Marcel Nonnen-
1259 macher, and Jakob H Macke. Flexible statistical inference for mechanistic models of neural
1260 dynamics. In *Advances in Neural Information Processing Systems*, pages 1289–1299, 2017.
- 1261 [89] George Papamakarios, David Sterratt, and Iain Murray. Sequential neural likelihood: Fast
1262 likelihood-free inference with autoregressive flows. In *The 22nd International Conference on*
1263 *Artificial Intelligence and Statistics*, pages 837–848. PMLR, 2019.
- 1264 [90] Joeri Hermans, Volodimir Begy, and Gilles Louppe. Likelihood-free mcmc with amortized
1265 approximate ratio estimators. In *International Conference on Machine Learning*, pages 4239–
1266 4248. PMLR, 2020.
- 1267 [91] Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and
1268 variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.
- 1269 [92] Sean R Bittner and John P Cunningham. Approximating exponential family models (not
1270 single distributions) with a two-network architecture. *arXiv preprint arXiv:1903.07515*, 2019.
- 1271 [93] Johan Karlsson, Milena Anguelova, and Mats Jirstrand. An efficient method for structural
1272 identifiability analysis of large dynamic systems. *IFAC Proceedings Volumes*, 45(16):941–946,
1273 2012.

- 1274 [94] Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary
1275 differential equations. In *Advances in neural information processing systems*, pages 6571–6583,
1276 2018.
- 1277 [95] Xuechen Li, Ting-Kam Leonard Wong, Ricky TQ Chen, and David Duvenaud. Scalable
1278 gradients for stochastic differential equations. *arXiv preprint arXiv:2001.01328*, 2020.
- 1279 [96] Andreas Raue, Clemens Kreutz, Thomas Maiwald, Julie Bachmann, Marcel Schilling, Ursula
1280 Klingmüller, and Jens Timmer. Structural and practical identifiability analysis of partially
1281 observed dynamical models by exploiting the profile likelihood. *Bioinformatics*, 25(15):1923–
1282 1929, 2009.
- 1283 [97] Dhruva V Raman, James Anderson, and Antonis Papachristodoulou. Delineating parameter
1284 unidentifiabilities in complex models. *Physical Review E*, 95(3):032314, 2017.
- 1285 [98] Maria Pia Saccomani, Stefania Audoly, and Leontina D’Angiò. Parameter identifiability of
1286 nonlinear systems: the role of initial conditions. *Automatica*, 39(4):619–632, 2003.
- 1287 [99] George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Bal-
1288 aji Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *arXiv
1289 preprint arXiv:1912.02762*, 2019.
- 1290 [100] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolu-
1291 tions. In *Advances in neural information processing systems*, pages 10215–10224, 2018.
- 1292 [101] Durk P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling.
1293 Improved variational inference with inverse autoregressive flow. *Advances in neural informa-
1294 tion processing systems*, 29:4743–4751, 2016.
- 1295 [102] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *Inter-
1296 national Conference on Learning Representations*, 2015.
- 1297 [103] David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for
1298 statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.

1299 **5 Methods**

1300 **5.1 Emergent property inference (EPI)**

1301 Determining the combinations of model parameters that can produce observed data or a desired
1302 output is a key part of scientific practice. Solving inverse problems is especially important in
1303 neuroscience, since we require complex models to describe the complex phenomena of neural com-
1304 putations. While much machine learning research has focused on how to find latent structure
1305 in large-scale neural datasets, less has focused on inverting theoretical circuit models conditioned
1306 upon the emergent phenomena they produce. Here, we introduce a novel method for statistical
1307 inference, which finds distributions of parameter solutions that only produce the desired emer-
1308 gent property. This method seamlessly handles neural circuit models with stochastic nonlinear
1309 dynamical generative processes, which are predominant in theoretical neuroscience.

1310 Consider model parameterization \mathbf{z} , which is a collection of scientifically interesting variables that
1311 govern the complex simulation of data \mathbf{x} . For example (see Section 3.1 Motivating emergent property
1312 inference of theoretical models subsection 3.1), \mathbf{z} may be the electrical conductance parameters of
1313 an STG subcircuit, and \mathbf{x} the evolving membrane potentials of the five neurons. In terms of
1314 statistical modeling, this circuit model has an intractable likelihood $p(\mathbf{x} | \mathbf{z})$, which is predicated
1315 by the stochastic differential equations that define the model. Even so, we do not scientifically
1316 reason about how \mathbf{z} governs all of \mathbf{x} , but rather specific phenomena that are a function of the
1317 data $f(\mathbf{x}; \mathbf{z})$. In the STG example, $f(\mathbf{x}; \mathbf{z})$ measures hub neuron frequency from the evolution of
1318 \mathbf{x} governed by \mathbf{z} . With EPI, we learn distributions of \mathbf{z} that results in an average and variance of
1319 $f(\mathbf{x}; \mathbf{z})$, denoted $\boldsymbol{\mu}$ and σ^2 . We refer to the collection of these statistical moments as an emergent
1320 property. Such emergent properties \mathcal{X} are defined through choice of $f(\mathbf{x}; \mathbf{z})$ (which may be one or
1321 multiple statistics), $\boldsymbol{\mu}$, and σ^2

$$\mathcal{X} : \mathbb{E}_{\mathbf{z}, \mathbf{x}} [f(\mathbf{x}; \mathbf{z})] = \boldsymbol{\mu}, \text{Var}_{\mathbf{z}, \mathbf{x}} [f(\mathbf{x}; \mathbf{z})] = \sigma^2. \quad (12)$$

1322 Precisely, the emergent property statistics $f(\mathbf{x}; \mathbf{z})$ must have means $\boldsymbol{\mu}$ and variances σ^2 over the
1323 EPI distribution of parameters and stochasticity of the data given the parameters.

1324 In EPI, deep probability distributions are used as posterior approximations $q_{\boldsymbol{\theta}}(\mathbf{z} | \mathcal{X})$. In deep
1325 probability distributions, a simple random variable $\mathbf{z}_0 \sim q_0(\mathbf{z}_0)$ is mapped deterministically via a
1326 sequence of deep neural network layers (g_1, \dots, g_l) parameterized by weights and biases $\boldsymbol{\theta}$ to the

1327 support of the distribution of interest:

$$\mathbf{z} = g_{\theta}(\mathbf{z}_0) = g_l(\dots g_1(\mathbf{z}_0)) \sim q_{\theta}(\mathbf{z}). \quad (13)$$

1328 Such deep probability distributions embed the posterior distribution in a deep network. Once
1329 optimized, this deep network representation has remarkably useful properties: immediate posterior
1330 sampling, and immediate probability, gradient, and Hessian evaluation at any parameter choice.

1331 Given a choice of model $p(\mathbf{x} | \mathbf{z})$ and emergent property of interest \mathcal{X} , $q_{\theta}(\mathbf{z})$ is optimized via
1332 the neural network parameters θ to find a maximally entropic distribution q_{θ}^* within the deep
1333 variational family \mathcal{Q} producing the emergent property \mathcal{X} :

$$q_{\theta}(\mathbf{z} | \mathcal{X}) = q_{\theta}^*(\mathbf{z}) = \underset{\theta \in Q}{\operatorname{argmax}} H(q_{\theta}(\mathbf{z})) \quad (14)$$
$$\text{s.t. } \mathcal{X} : \mathbb{E}_{\mathbf{z}, \mathbf{x}} [f(\mathbf{x}; \mathbf{z})] = \boldsymbol{\mu}, \operatorname{Var}_{\mathbf{z}, \mathbf{x}} [f(\mathbf{x}; \mathbf{z})] = \boldsymbol{\sigma}^2.$$

1334 Entropy is chosen as the normative selection principle, since we want the posterior to only contain
1335 structure predicated by the emergent property [58, 59]. This choice of selection principle is also that
1336 of standard Bayesian inference, and we derive an exact relation between EPI and variational infer-
1337 ence (see Section 5.1.5Maximum entropy distributions and exponential familiessubsubsection.5.1.5).
1338 However, a key difference is that variational inference and other Bayesian methods do not constrain
1339 the predictions of their inferred posteriors. This optimization is executed using the algorithm of
1340 Maximum Entropy Flow Networks (MEFNs) [52].

1341 In the remainder of Section 5.1Emergent property inference (EPI)subsection.5.1, we will explain
1342 the finer details and motivation of the EPI method. First, we explain related approaches and what
1343 EPI introduces to this domain (Section 5.1.1Related approachessubsubsection.5.1.1). Second, we
1344 describe the special class of deep probability distributions used in EPI called normalizing flows
1345 (Section 5.1.2Normalizing flowssubsubsection.5.1.2). Next, we explain the constrained optimiza-
1346 tion technique used to solve Equation 7Emergent property inference (EPI)equation.5.7 (Section
1347 5.1.3Augmented Lagrangian optimizationsubsubsection.5.1.3). Then, we demonstrate the details
1348 of this optimization in a toy example (Section 5.1.4Example: 2D LDSsubsubsection.5.1.4). Finally,
1349 we establish the known relationship between maximum entropy distributions and exponential fam-
1350 ilies (Section 5.1.5Maximum entropy distributions and exponential familiessubsubsection.5.1.5),
1351 which is used to explain the relation between EPI and variational inference (Section 5.1.6EPI as
1352 variational inferencesubsubsection.5.1.6).

1353 **5.1.1 Related approaches**

1354 When Bayesian inference problems lack conjugacy, scientists use approximate inference methods
1355 like variational inference (VI) [46] and Markov chain Monte Carlo (MCMC) [45, 44]. After optimi-
1356 zation, variational methods return a parameterized posterior distribution, which we can analyze.
1357 Also, the variational approximating distribution class is often chosen such that it permits fast
1358 sampling. In contrast MCMC methods only produce samples from the approximated posterior dis-
1359 tribution. No parameterized distribution is estimated, and additional samples are always generated
1360 with the same sampling complexity. Inference in models defined by systems of differential has been
1361 demonstrated with MCMC [81], although this approach requires tractable likelihoods. Advance-
1362 ments have leveraged structure in stochastic differential equation models to improve likelihood
1363 approximations, thus expanding the domain of applicable models [82].

1364 Likelihood-free (or “simulation-based”) inference (LFI) [83] is model parameter inference in the
1365 absence of a tractable likelihood function. The most prevalent approach to LFI is approximate
1366 Bayesian computation [42], in which satisfactory parameter samples are kept from random prior
1367 sampling according to a rejection heuristic. The obtained set of parameters do not have a probabili-
1368 ties, and further insight about the model must be gained from examination of the parameter set and
1369 their generated activity. Methodological advances to ABC methods have come through the use of
1370 Markov chain Monte Carlo (MCMC-ABC) [43] and sequential Monte Carlo (SMC-ABC) [31] sam-
1371 pling techniques. SMC-ABC is considered state-of-the-art ABC, yet this approach still struggles to
1372 scale in dimensionality (cf. Fig. 4**A**. Recurrent neural network. **B**. EPI scales with z to high dimensions.
1373 Convergence definitions: EPI (blue) - satisfies all moment constraints, SNPE (orange)- produces at least
1374 $2/n_{\text{train}}$ parameter samples are in the bounds of emergent property (mean ± 0.5), and SMC-ABC (red) - 100
1375 particles with $\epsilon < 0.5$ are produced. **C**. Posterior predictive distributions of EPI (blue), SNPE (orange), and
1376 SMC-ABC (red). Gray star indicates emergent property mean, and gray dashed lines indicate two standard
1377 deviations corresponding to the variance constraint. For $N \leq 6$ where SMC-ABC converges, samples are
1378 not diverse (path degeneracies). For $N \geq 25$, SNPE does not produce a posterior approximation yielding
1379 parameters with simulations near x_0 . **D**. Simulations of network parameters resulting from each method
1380 ($\tau = 100ms$). Each trace corresponds to simulation of one z . (Below) Ratio of obtained samples producing
1381 stable amplificationfigure.4). Furthermore, once a parameter set has been obtained by SMC-ABC
1382 from a finite set of particles, the SMC-ABC algorithm must be run again with a new population
1383 of initialized particles to obtain additional samples.
1384 For scientific model analysis, we seek a posterior distribution exhibiting the properties of a well-

1385 chosen variational approximation: a parametric form conferring analytic calculations, and trivial
1386 sampling time. For this reason, ABC and MCMC techniques are unattractive, since they only
1387 produce a set of parameter samples and have unchanging sampling rate. EPI executes likelihood-
1388 free inference using the MEFN [52] algorithm using a deep variational posterior approximation.
1389 The deep neural network of EPI defines the parametric form of the posterior approximation. Fur-
1390 thermore, the EPI distribution is constrained to produce an emergent property. In other words,
1391 the summary statistics of the posterior predictive distribution are fixed to have certain first and
1392 second moments. EPI optimization is enabled using stochastic gradient techniques in the spirit
1393 of likelihood-free variational inference [49]. The analytic relationship between EPI and variational
1394 inference is explained in Secton 5.1.6EPI as variational inferencesubsubsection.5.1.6.

1395 We note that, during our preparation and early presentation of this work [84, 85], another work
1396 has arisen with broadly similar goals: bringing statistical inference to mechanistic models of neural
1397 circuits ([86, 87, 75]). We are encouraged by this general problem being recognized by others in the
1398 community, and we emphasize that these works offer complementary neuroscientific contributions
1399 (different theoretical models of focus) and use different technical methodologies (ours is built on
1400 our prior work [52], theirs similarly [88]).

1401 The method EPI differs from SNPE in some key ways. SNPE belongs to a “sequential” class of
1402 recently developed LFI methods in which two neural networks are used for posterior inference.
1403 This first neural network is a normalizing flow used to estimate the posterior $p(\mathbf{z} | \mathbf{x})$ (SNPE)
1404 or the likelihood $p(\mathbf{x} | \mathbf{z})$ (sequential neural likelihood (SNL [89])). A recent advance uses an
1405 unconstrained neural network to estimate the likelihood ratio (sequential neural ratio estimation
1406 (SNRE [90])). In SNL and SNRE, MCMC sampling techniques are used to obtain samples from
1407 the approximated posterior. This contrasts with EPI and SNPE, which afford a normalizing flow
1408 approximation to the posterior, which facilitates immediate measurements of sample probability,
1409 gradient, or Hessian for system analysis. The second neural network in this sequential class of
1410 methods is the amortizer. This network maps data \mathbf{x} (or statistics $f(\mathbf{x}; \mathbf{z})$ or model parameters \mathbf{z})
1411 to the weights and biases of the first neural network. These methods are optimized on a conditional
1412 density (or ratio) estimation objective on a sequentially adapting finite sample-based approximation
1413 to the posterior.

1414 The approximating fidelity of the first neural network in sequential approaches is optimized to
1415 generalize across the entire distribution it is conditioned upon. This optimization towards gen-
1416 eralization of sequential methods can reduce the accuracy at the singular posterior of interest.

Whereas in EPI, the entire expressivity of the normalizing flow is dedicated to learning a single distribution as well as possible. While amortization is not possible in EPI parameterized by the mean parameter μ (due to the inverse mapping problem [91]), we have shown this two-network amortization approach to be effective in exponential family distributions defined by their natural parameterization [92].

Structural identifiability analysis involves the measurement of sensitivity and unidentifiabilities in natural models. Around a point, one can measure the Jacobian. One approach that scales well is EAR [93]. A popular efficient approach for systems of ODEs has been neural ODE adjoint [94] and its stochastic adaptation [95]. Casting identifiability as a statistical estimation problem, the profile likelihood can assess via iterated optimization while holding parameters fixed [96]. An exciting recent method is capable of recovering the functional form of such unidentifiabilities away from a point by following degenerate dimensions of the fisher information matrix [97]. Global structural non-identifiabilities can be found for models with polynomial or rational dynamics equations using DAISY [98]. With EPI, we have all the benefits given by a statistical inference method plus the ability to query the gradient or Hessian of the inferred distribution at any chosen parameter value.

5.1.2 Normalizing flows

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

$$p(\mathbf{z}_i) = p(g_i^{-1}(\mathbf{z}_i)) \left| \det \frac{\partial g_i^{-1}(\mathbf{z}_i)}{\partial \mathbf{z}_i} \right| = p(\mathbf{z}_{i-1}) \left| \det \frac{\partial g_i(\mathbf{z}_{i-1})}{\partial \mathbf{z}_{i-1}} \right|^{-1}. \quad (15)$$

However, this computation has cubic complexity in dimensionality for fully connected layers. By restricting our layers to normalizing flows [47, 99] – bijective functions with fast log determinant Jacobian computations, which confer a fast calculation of the sample log probability. Fast log probability calculation confers efficient optimization of the maximum entropy objective (see Section 5.1.3 Augmented Lagrangian optimizationsubsubsection 5.1.3). We use the Real NVP [50] normalizing flow class, because its coupling architecture confers both fast sampling (forward) and fast log probability evaluation (backward). Fast probability evaluation in turn facilitates fast gradient and Hessian evaluation of log probability throughout parameter space. Glow permutations were used in between coupling stages [100]. This is in contrast to autoregressive architectures [51, 101], in

1446 which only forward or backward passes are efficient. In this work, normalizing flows are used as
 1447 flexible posterior approximations $q_{\theta}(\mathbf{z})$ having weights and biases $\boldsymbol{\theta}$. We specify the architecture
 1448 used in each application by the number of Real-NVP affine coupling stages, and the number of
 1449 neural network layers and units per layer of the conditioning functions.

1450 5.1.3 Augmented Lagrangian optimization

1451 To optimize $q_{\theta}(\mathbf{z})$ in Equation 7Emergent property inference (EPI)equation.5.7, the constrained
 1452 maximum entropy optimization is executed using the augmented Lagrangian method. The following
 1453 objective is minimized:

$$L(\boldsymbol{\theta}; \boldsymbol{\eta}_{\text{opt}}, c) = -H(q_{\theta}) + \boldsymbol{\eta}_{\text{opt}}^T R(\boldsymbol{\theta}) + \frac{c}{2} \|R(\boldsymbol{\theta})\|^2 \quad (16)$$

1454 where average constraint violations $R(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{z} \sim q_{\theta}} [\mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}|\mathbf{z})} [T(\mathbf{x}; \mathbf{z}) - \boldsymbol{\mu}_{\text{opt}}]]$, $\boldsymbol{\eta}_{\text{opt}} \in \mathbb{R}^m$ are the
 1455 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.
 1456 The sufficient statistics $T(\mathbf{x}; \mathbf{z})$ and mean parameter $\boldsymbol{\mu}_{\text{opt}}$ are determined by the means $\boldsymbol{\mu}$ and vari-
 1457 ances $\boldsymbol{\sigma}^2$ of emergent property statistics $f(\mathbf{x}; \mathbf{z})$ defined in Equation 7Emergent property inference
 1458 (EPI)equation.5.7. Specifically, $T(\mathbf{x}; \mathbf{z})$ is a concatenation of the first and second moments, $\boldsymbol{\mu}_{\text{opt}}$
 1459 is a concatenation of $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}^2$ (see section 5.1.5Maximum entropy distributions and exponential
 1460 familiessubsubsection.5.1.5), and the Lagrange multipliers are closely related to the natural param-
 1461 eters $\boldsymbol{\eta}$ of exponential families (see Section 5.1.6EPI as variational inferencesubsubsection.5.1.6).
 1462 Weights and biases $\boldsymbol{\theta}$ of the deep probability distribution are optimized according to Equation ??
 1463 using the Adam optimizer with learning rate 10^{-3} [102].

1464 To take gradients with respect to the entropy $H(q_{\theta}(\mathbf{z}))$, it can be expressed using the reparam-
 1465 eterization trick as an expectation of the negative log density of parameter samples \mathbf{z} over the
 1466 randomness in the parameterless initial distribution $q_0(\mathbf{z}_0)$:

$$H(q_{\theta}(\mathbf{z})) = \int -q_{\theta}(\mathbf{z}) \log(q_{\theta}(\mathbf{z})) d\mathbf{z} = \mathbb{E}_{\mathbf{z} \sim q_{\theta}} [-\log(q_{\theta}(\mathbf{z}))] = \mathbb{E}_{\mathbf{z}_0 \sim q_0} [-\log(q_{\theta}(g_{\theta}(\mathbf{z}_0)))]. \quad (17)$$

1467 Thus, the gradient of the entropy of the deep probability distribution can be estimated as an
 1468 average with respect to the base distribution \mathbf{z}_0 :

$$\nabla_{\boldsymbol{\theta}} H(q_{\theta}(\mathbf{z})) = \mathbb{E}_{\mathbf{z}_0 \sim q_0} [-\nabla_{\boldsymbol{\theta}} \log(q_{\theta}(g_{\theta}(\mathbf{z}_0)))] . \quad (18)$$

1469 The lagrangian parameters $\boldsymbol{\eta}_{\text{opt}}$ are initialized to zero and adapted following each augmented
 1470 Lagrangian epoch, which is a period of optimization with fixed $(\boldsymbol{\eta}_{\text{opt}}, c)$ for a given number of

1471 stochastic optimization iterations. A low value of c is used initially, and conditionally increased
1472 after each epoch based on constraint error reduction. The penalty coefficient is updated based
1473 on the result of a hypothesis test regarding the reduction in constraint violation. The p-value of
1474 $\mathbb{E}[|R(\boldsymbol{\theta}_{k+1})|] > \gamma \mathbb{E}[|R(\boldsymbol{\theta}_k)|]$ is computed, and c_{k+1} is updated to βc_k with probability $1 - p$. The
1475 other update rule is $\boldsymbol{\eta}_{\text{opt},k+1} = \boldsymbol{\eta}_{\text{opt},k} + c_k \frac{1}{n} \sum_{i=1}^n (T(\mathbf{x}^{(i)}) - \boldsymbol{\mu}_{\text{opt}})$ given a batch size n . Throughout
1476 the study, $\gamma = 0.25$, while β was chosen to be either 2 or 4. The batch size of EPI also varied
1477 according to application.

1478 The intention is that c and $\boldsymbol{\eta}_{\text{opt}}$ start at values encouraging entropic growth early in optimization.
1479 With each training epoch in which the update rule for c is invoked by unsatisfactory constraint
1480 error reduction, the constraint satisfaction terms are increasingly weighted, resulting in a decreased
1481 entropy. This encourages the discovery of suitable regions of parameter space, and the subsequent
1482 refinement of the distribution to produce the emergent property (see example in Section 5.1.4Ex-
1483 ample: 2D LDSsubsubsection.5.1.4). The momentum parameters of the Adam optimizer are reset
1484 at the end of each augmented Lagrangian epoch.

1485 Rather than starting optimization from some $\boldsymbol{\theta}$ drawn from a randomized distribution, we found
1486 that initializing $q_{\boldsymbol{\theta}}(\mathbf{z})$ to approximate an isotropic Gaussian distribution conferred more stable, con-
1487 sistent optimization. The parameters of the Gaussian initialization were chosen on an application-
1488 specific basis. Throughout the study, we chose isotropic Gaussian initializations with mean $\boldsymbol{\mu}_{\text{init}}$ at
1489 the center of the distribution support and some standard deviation σ_{init} , except for one case,
1490 where an initialization informed by random search was used (see Section 5.2.1Stomatogastric
1491 ganglionssubsubsection.5.2.1).

1492 To assess whether the EPI distribution $q_{\boldsymbol{\theta}}(\mathbf{z})$ produces the emergent property, we assess whether
1493 each individual constraint on the means and variances of $f(\mathbf{x}; \mathbf{z})$ is satisfied. We consider the EPI
1494 to have converged when a null hypothesis test of constraint violations $R(\boldsymbol{\theta})_i$ being zero is accepted
1495 for all constraints $i \in \{1, \dots, m\}$ at a significance threshold $\alpha = 0.05$. This significance threshold is
1496 adjusted through Bonferroni correction according to the number of constraints m . The p-values for
1497 each constraint are calculated according to a two-tailed nonparametric test, where 200 estimations
1498 of the sample mean $R(\boldsymbol{\theta})^i$ are made using N_{test} samples of $\mathbf{z} \sim q_{\boldsymbol{\theta}}(\mathbf{z})$ at the end of the augmented
1499 Lagrangian epoch.

1500 When assessing the suitability of EPI for a particular modeling question, there are some important
1501 technical considerations. First and foremost, as in any optimization problem, the defined emergent
1502 property should always be appropriately conditioned (constraints should not have wildly different

units). Furthermore, if the program is underconstrained (not enough constraints), the distribution grows (in entropy) unstably unless mapped to a finite support. If overconstrained, there is no parameter set producing the emergent property, and EPI optimization will fail (appropriately). Next, one should consider the computational cost of the gradient calculations. In the best circumstance, there is a simple, closed form expression (e.g. Section 5.2.5Rank-2 RNNsubsubsection.5.2.5) for the emergent property statistic given the model parameters. On the other end of the spectrum, many forward simulation iterations may be required before a high quality measurement of the emergent property statistic is available (e.g. Section 5.2.1Stomatogastric ganglionsubsubsection.5.2.1). In such cases, backpropagating gradients through the SDE evolution will be expensive.

5.1.4 Example: 2D LDS

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

$$\tau \frac{d\mathbf{x}}{dt} = A\mathbf{x} \quad (19)$$

with

$$A = \begin{bmatrix} a_1 & a_2 \\ a_3 & a_4 \end{bmatrix}. \quad (20)$$

To run EPI with the dynamics matrix elements as the free parameters $\mathbf{z} = [a_1, a_2, a_3, a_4]$ (fixing $\tau = 1$), the emergent property statistics $T(\mathbf{x})$ were chosen to contain the first and second moments of the oscillatory frequency, $\frac{\text{imag}(\lambda_1)}{2\pi}$, and the growth/decay factor, $\text{real}(\lambda_1)$, of the oscillating system. λ_1 is the eigenvalue of greatest real part when the imaginary component is zero, and alternatively of positive imaginary component when the eigenvalues are complex conjugate pairs. To learn the distribution of real entries of A that produce a band of oscillating systems around 1Hz, we formalized this emergent property as $\text{real}(\lambda_1)$ having mean zero with variance 0.25^2 , and the oscillation frequency $2\pi\text{imag}(\lambda_1)$ having mean $\omega = 1$ Hz with variance $(0.1\text{Hz})^2$:

$$\mathbb{E}[T(\mathbf{x})] \triangleq \mathbb{E} \begin{bmatrix} \text{real}(\lambda_1) \\ \text{imag}(\lambda_1) \\ (\text{real}(\lambda_1) - 0)^2 \\ (\text{imag}(\lambda_1) - 2\pi\omega)^2 \end{bmatrix} = \begin{bmatrix} 0.0 \\ 2\pi\omega \\ 0.25^2 \\ (2\pi 0.1)^2 \end{bmatrix} \triangleq \boldsymbol{\mu}. \quad (21)$$

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

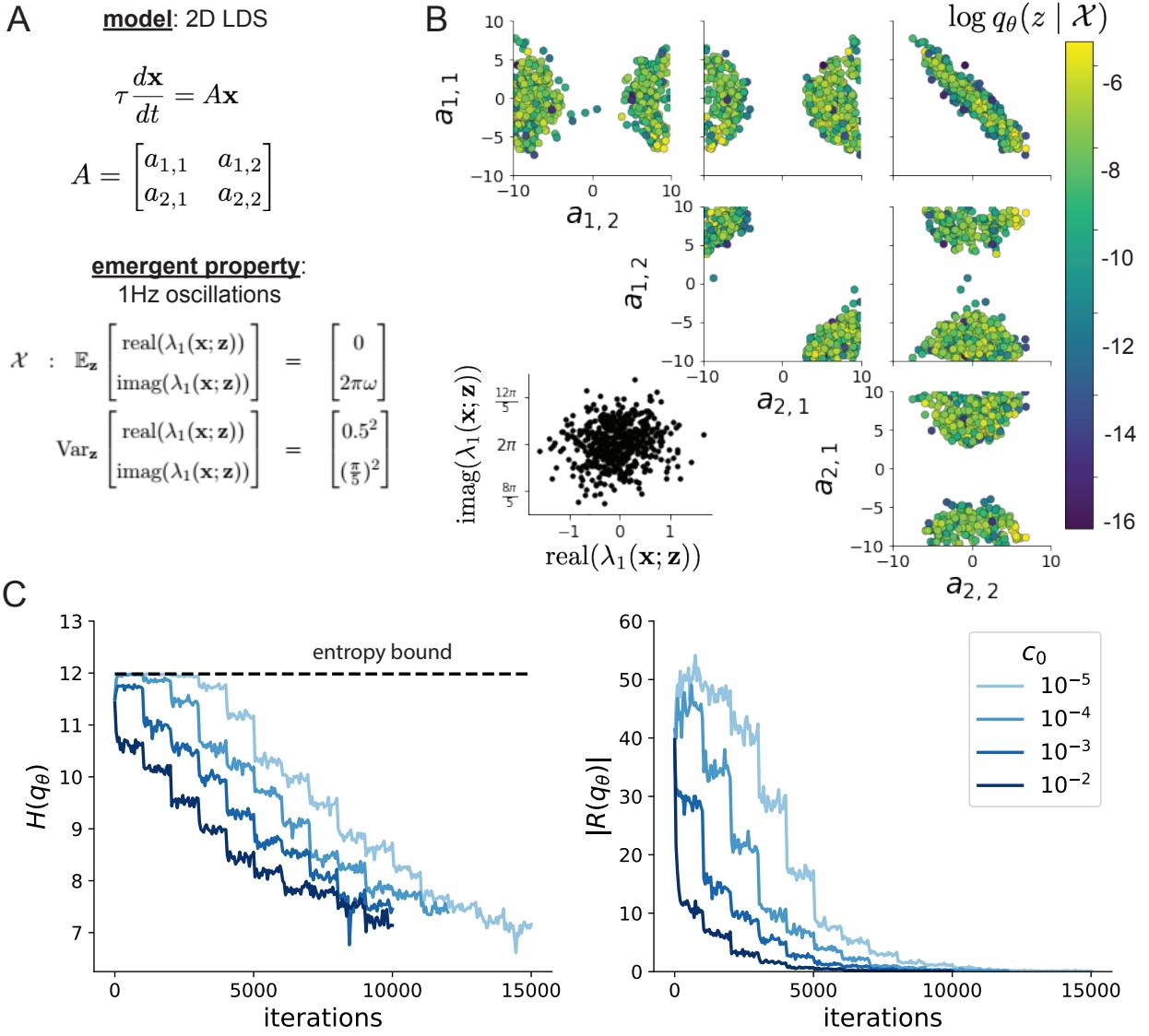


Figure 5: (LDS1): A. Two-dimensional linear dynamical system model, where real entries of the dynamics matrix A are the parameters. B. The EPI distribution for a two-dimensional linear dynamical system with $\tau = 1$ that produces an average of 1Hz oscillations with some small amount of variance. Dashed 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.

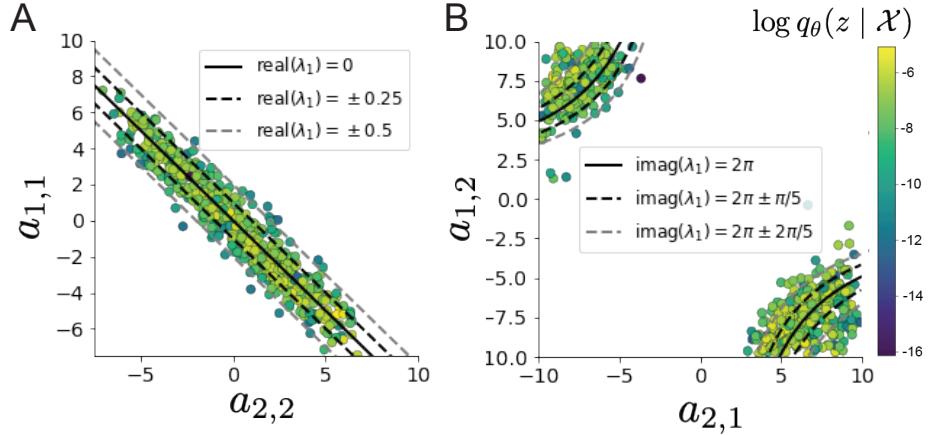


Figure 6: (LDS2): A. Probability contours in the a_1 - a_4 plane were derived from the relationship to emergent property statistic of growth/decay factor $\text{real}(\lambda_1)$. B. Probability contours in the a_2 - a_3 plane were derived from the emergent property statistic of oscillation frequency $2\pi\text{imag}(\lambda_1)$.

1527 the quadratic formula:

$$\lambda = \frac{\left(\frac{a_1+a_4}{\tau}\right) \pm \sqrt{\left(\frac{a_1+a_4}{\tau}\right)^2 + 4\left(\frac{a_2a_3-a_1a_4}{\tau}\right)}}{2}. \quad (22)$$

1528 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
 1529 does not require simulation, we cannot derive the distribution q_θ^* directly. This fact is due to the
 1530 formally hard problem of the backward mapping: finding the natural parameters η from the mean
 1531 parameters μ of an exponential family distribution [91]. Instead, we used EPI to approximate this
 1532 distribution (Fig. S1B). We used a real-NVP normalizing flow architecture with four masks, two
 1533 neural network layers of 15 units per mask, with batch normalization momentum 0.99, mapped
 1534 onto a support of $z_i \in [-10, 10]$. (see Section 5.1.2Normalizing flowssubsubsection.5.1.2).

1535 Even this relatively simple system has nontrivial (though intuitively sensible) structure in the
 1536 parameter distribution. To validate our method, we analytically derived the contours of the prob-
 1537 ability density from the emergent property statistics and values. In the a_1 - a_4 plane, the black
 1538 line at $\text{real}(\lambda_1) = \frac{a_1+a_4}{2} = 0$, dotted black line at the standard deviation $\text{real}(\lambda_1) = \frac{a_1+a_4}{2} \pm 0.25$,
 1539 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
 1540 of probability density of the samples (Fig. S2A). The distribution precisely reflects the desired
 1541 statistical constraints and model degeneracy in the sum of a_1 and a_4 . Intuitively, the parameters
 1542 equivalent with respect to emergent property statistic $\text{real}(\lambda_1)$ have similar log densities.

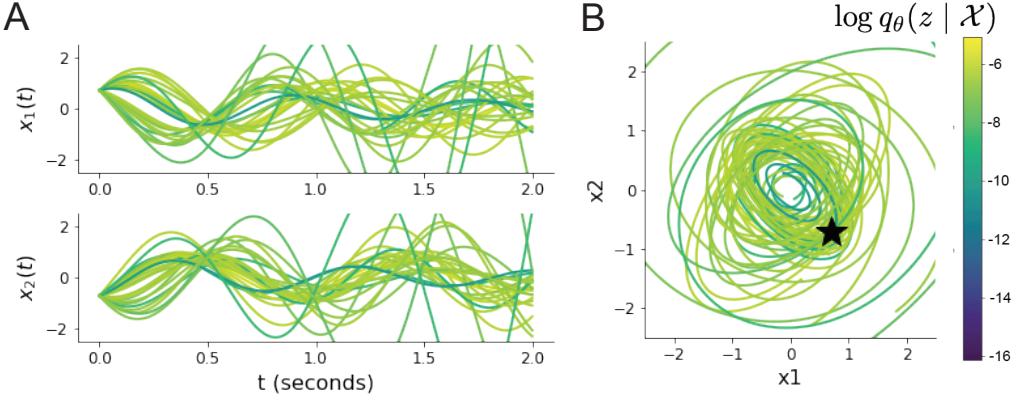


Figure 7: (LDS3): Sampled dynamical systems $\mathbf{z} \sim q_{\theta}(\mathbf{z})$ and their simulated activity from $\mathbf{x}(0) = [\frac{\sqrt{2}}{2}, -\frac{\sqrt{2}}{2}]$ colored by log probability. A. Each dimension of the simulated trajectories throughout time. B. The simulated trajectories in phase space.

1543 To explain the bimodality of the EPI distribution, we examined the imaginary component of λ_1 .

1544 When $\text{real}(\lambda_1) = \frac{a_1 + a_4}{2} = 0$, we have

$$\text{imag}(\lambda_1) = \begin{cases} \sqrt{\frac{a_1 a_4 - a_2 a_3}{\tau}}, & \text{if } a_1 a_4 < a_2 a_3 \\ 0 & \text{otherwise} \end{cases}. \quad (23)$$

1545 When $\tau = 1$ and $a_1 a_4 > a_2 a_3$ (center of distribution above), we have the following equation for the
1546 other two dimensions:

$$\text{imag}(\lambda_1)^2 = a_1 a_4 - a_2 a_3 \quad (24)$$

1547 Since we constrained $\mathbb{E}_{\mathbf{z} \sim q_{\theta}} [\text{imag}(\lambda)] = 2\pi$ (with $\omega = 1$), we can plot contours of the equation
1548 $\text{imag}(\lambda_1)^2 = a_1 a_4 - a_2 a_3 = (2\pi)^2$ for various $a_1 a_4$ (Fig. S2B). With $\sigma_{1,4} = \mathbb{E}_{\mathbf{z} \sim q_{\theta}} (|a_1 a_4 - E_{q_{\theta}}[a_1 a_4]|)$,
1549 we show the contours as $a_1 a_4 = 0$ (black), $a_1 a_4 = -\sigma_{1,4}$ (black dotted), and $a_1 a_4 = -2\sigma_{1,4}$ (grey
1550 dotted). This validates the curved structure of the inferred distribution learned through EPI. We
1551 took steps in negative standard deviation of $a_1 a_4$ (dotted and gray lines), since there are few positive
1552 values $a_1 a_4$ in the learned distribution. Subtler combinations of model and emergent property will
1553 have more complexity, further motivating the use of EPI for understanding these systems. As we
1554 expect, the distribution results in samples of two-dimensional linear systems oscillating near 1Hz
1555 (Fig. S3).

1556 **5.1.5 Maximum entropy distributions and exponential families**

1557 Maximum entropy distributions have a fundamental link to exponential family distributions. A
 1558 maximum entropy distribution of form:

$$p^*(\mathbf{z}) = \underset{p \in \mathcal{P}}{\operatorname{argmax}} H(p(\mathbf{z})) \quad (25)$$

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

1559 will have probability density in the exponential family:

$$p^*(\mathbf{z}) \propto \exp(\boldsymbol{\eta}^\top T(\mathbf{z})). \quad (26)$$

1560 The mappings between the mean parameterization $\boldsymbol{\mu}_{\text{opt}}$ and the natural parameterization $\boldsymbol{\eta}$ are
 1561 formally hard to identify [91].

1562 In EPI, emergent properties are defined as statistics having a fixed mean and variance as in Equation
 1563 2A deep generative modeling approach to emergent property inferenceequation.3.2

$$\mathbb{E}_{\mathbf{z}, \mathbf{x}}[f(\mathbf{x}; \mathbf{z})] = \boldsymbol{\mu}, \operatorname{Var}_{\mathbf{z}, \mathbf{x}}[f(\mathbf{x}; \mathbf{z})] = \sigma^2. \quad (27)$$

1564 The variance constraint is a second moment constraint on $f(\mathbf{x}; \mathbf{z})$

$$\operatorname{Var}_{\mathbf{z}, \mathbf{x}}[f(\mathbf{x}; \mathbf{z})] = \mathbb{E}_{\mathbf{z}, \mathbf{x}}[(f(\mathbf{x}; \mathbf{z}) - \boldsymbol{\mu})^2] \quad (28)$$

1565 As a general maximum entropy distribution (Equation 17Maximum entropy distributions and ex-
 1566 ponential familiesequation.5.17), the sufficient statistics vector contains both first and second order
 1567 moments of $f(\mathbf{x}; \mathbf{z})$

$$T(\mathbf{x}; \mathbf{z}) = \begin{bmatrix} f(\mathbf{x}; \mathbf{z}) \\ (f(\mathbf{x}; \mathbf{z}) - \boldsymbol{\mu})^2 \end{bmatrix}, \quad (29)$$

1568 which are constrained to the chosen means and variances

$$\boldsymbol{\mu}_{\text{opt}} = \begin{bmatrix} \boldsymbol{\mu} \\ \sigma^2 \end{bmatrix}. \quad (30)$$

1569 **5.1.6 EPI as variational inference**

1570 In Bayesian inference a prior belief about model parameters \mathbf{z} is stated in a prior distribution $p(\mathbf{z})$,
 1571 and the statistical model capturing the effect of \mathbf{z} on observed data points \mathbf{x} is formalized in the

1572 likelihood distribution $p(\mathbf{x} \mid \mathbf{z})$. In Bayesian inference, we obtain a posterior distribution $p(z \mid \mathbf{x})$,
 1573 which captures how the data inform our knowledge of model parameters using Bayes' rule:

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

1574 The posterior distribution is analytically available when the prior is conjugate with the likelihood.
 1575 However, conjugacy is rare in practice, and alternative methods, such as variational inference [103],
 1576 are utilized.

1577 In variational inference, a posterior approximation $q_{\boldsymbol{\theta}}^*$ is chosen from within some variational family
 1578 \mathcal{Q}

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

1579 The KL divergence can be written in terms of entropy of the variational approximation:

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

1580

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

1581 Since the marginal distribution of the data $p(\mathbf{x})$ (or “evidence”) is independent of $\boldsymbol{\theta}$, variational
 1582 inference is executed by optimizing the remaining expression. This is usually framed as maximizing
 1583 the evidence lower bound (ELBO)

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

1584 Now, consider the setting where we have chosen a uniform prior, and stipulate a mean-field gaussian
 1585 likelihood on a chosen statistic of the data $f(\mathbf{x}; \mathbf{z})$

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

1586 where $\Sigma_f = \operatorname{diag}(\boldsymbol{\sigma}_f^2)$. The log likelihood is then proportional to a dot product of the natural
 1587 parameter of this mean-field gaussian distribution and the first and second moment statistics.

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

1588 where

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

1589

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

1590 The variational objective is then

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

1591 Comparing this to the Lagrangian objective (without augmentation) of EPI, we see they are the

1592 same

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

1593 where $T(\mathbf{x}; \mathbf{z})$ consists of the first and second moments of the emergent property statistic $f(\mathbf{x}; \mathbf{z})$
1594 (Equation ??). Thus, EPI is implicitly executing variational inference with a uniform prior and a
1595 mean-field gaussian likelihood on the emergent property statistics. The data \mathbf{x} used by this implicit
1596 variational inference program would be that generated by the adapting variational approximation
1597 $\mathbf{x} \sim p(\mathbf{x} | \mathbf{z})q_{\theta}(\mathbf{z})$, and the likelihood parameters $\boldsymbol{\eta}_f$ of EPI optimization epoch k are predicated
1598 by $\boldsymbol{\eta}_{\text{opt}, k}$. However, in EPI we have not specified a prior distribution, or collected data, which can
1599 inform us about model parameters. Instead we have a mathematical specification of an emergent
1600 property, which the model must produce, and a maximum entropy selection principle. Accordingly,
1601 we replace the notation of $p(\mathbf{z} | \mathbf{x})$ with $p(\mathbf{z} | \mathcal{X})$ conceptualizing an inferred distribution that obeys
1602 emergent property \mathcal{X} (see Section 5.1Emergent property inference (EPI)subsection.5.1).

1603 5.2 Theoretical models

1604 In this study, we used emergent property inference to examine several models relevant to theoretical
1605 neuroscience. Here, we provide the details of each model and the related analyses.

1606 5.2.1 Stomatogastric ganglion

1607 We analyze how the parameters $\mathbf{z} = [g_{\text{el}}, g_{\text{synA}}]$ govern the emergent phenomena of intermediate
1608 hub frequency in a model of the stomatogastric ganglion (STG) [57] shown in Figure 1Emergent
1609 property inference (EPI) in the stomatogastric ganglion. **A.** Conductance-based biophysical model of the
1610 STG subcircuit. In the STG model, jagged connections indicate electrical coupling having electrical con-
1611 ductance g_{el} . Other connections in the diagram are inhibitory synaptic projections having strength g_{synA}
1612 onto the hub neuron, and $g_{\text{synB}} = 5\text{nS}$ for mutual inhibitory connections. Parameters are represented by the
1613 vector \mathbf{z} and membrane potentials by the vector \mathbf{x} . The evolution of this model's activity $\mathbf{x}(t)$ is predicated

1614 by differential equations. **B.** Spiking frequency $\omega(\mathbf{x}; \mathbf{z})$ is an emergent property statistic. In this example,
 1615 spiking frequency is measured from simulated activity of the STG model at parameter choices of $g_{el} = 4.5\text{nS}$
 1616 and $g_{synA} = 3\text{nS}$. **C.** The emergent property of intermediate hub frequency, in which the hub neuron fires
 1617 at a rate between the fast and slow frequencies. This emergent property is defined by a mean and variance
 1618 on the emergent property statistic. Simulated activity traces are colored by log probability density of their
 1619 generating parameters in the EPI-inferred distribution (Panel E). **D.** For a choice of model and emergent
 1620 property, emergent property inference (EPI) learns a deep probability distribution of parameters \mathbf{z} . Deep
 1621 probability distributions map a simple random variable \mathbf{z}_0 through a deep neural network with weights and
 1622 biases $\boldsymbol{\theta}$ to parameters $\mathbf{z} = g_{\boldsymbol{\theta}}(\mathbf{z}_0)$. In EPI optimization, stochastic gradient steps in $\boldsymbol{\theta}$ are taken such that
 1623 entropy is maximized, and the emergent property \mathcal{X} is produced. The EPI posterior distribution is denoted
 1624 $q_{\boldsymbol{\theta}}(\mathbf{z} | \mathcal{X})$. **E.** The EPI posterior producing intermediate hub frequency. Samples are colored by log probabili-
 1625 ty density. Distribution contours of average hub neuron frequency from mean of .55 Hz are shown at levels
 1626 of .525, .53,575 Hz (dark to light gray away from mean). Eigenvectors of the Hessian at the mode of the
 1627 inferred distribution are indicated as \mathbf{v}_1 (solid) and \mathbf{v}_2 (dashed) with lengths scaled by the square root of the
 1628 absolute value of their eigenvalues. **F** Simulations from parameters in E. (Top) The predictive distribution of
 1629 the posterior obeys the emergent property. The black and gray dashed lines show the mean and two standard
 1630 deviations according the emergent property, respectively. (Bottom) Simulations at the starred parameter
 1631 valuesfigure.1A with activity $\mathbf{x} = [x_{f1}, x_{f2}, x_{hub}, x_{s1}, x_{s2}]$, using the same hyperparameter choices as
 1632 Gutierrez et al. Each neuron's membrane potential $x_{\alpha}(t)$ for $\alpha \in \{f1, f2, hub, s1, s2\}$ is the solution
 1633 of the following stochastic differential equation:

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

1634 The input current of each neuron is the sum of the leak, calcium, potassium, hyperpolarization,
 1635 electrical and synaptic currents as well as gaussian noise dB . Each current component is a function
 1636 of all membrane potentials and the conductance parameters \mathbf{z} .
 1637 The capacitance of the cell membrane was set to $C_m = 1\text{nF}$. Specifically, the currents are the
 1638 difference in the neuron's membrane potential and that current type's reversal potential multiplied
 1639 by a conductance:

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

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

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

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

1643

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

1644

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

1645 The reversal potentials were set to $V_{leak} = -40mV$, $V_{Ca} = 100mV$, $V_K = -80mV$, $V_{hyp} = -20mV$,
 1646 and $V_{syn} = -75mV$. The other conductance parameters were fixed to $g_{leak} = 1 \times 10^{-4}\mu S$. g_{Ca} ,
 1647 g_K , and g_{hyp} had different values based on fast, intermediate (hub) or slow neuron. The fast
 1648 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
 1649 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
 1650 slow conductances had values $g_{Ca} = 8.5 \times 10^{-3}$, $g_K = 1.5 \times 10^{-2}$, and $g_{hyp} = 1.0 \times 10^{-2}$.

1651 Furthermore, the Calcium, Potassium, and hyperpolarization channels have time-dependent gating
 1652 dynamics dependent on steady-state gating variables M_∞ , N_∞ and H_∞ , respectively:

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

1653

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

1654

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

1655

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

1656

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

1657

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

1658

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

1659 where we set $v_1 = 0mV$, $v_2 = 20mV$, $v_3 = 0mV$, $v_4 = 15mV$, $v_5 = 78.3mV$, $v_6 = 10.5mV$,
 1660 $v_7 = -42.2mV$, $v_8 = 87.3mV$, $v_9 = 5mV$, and $v_{th} = -25mV$.

1661 Finally, there is a synaptic gating variable as well:

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

1662 When the dynamic gating variables are considered, this is actually a 15-dimensional nonlinear
 1663 dynamical system. Gaussian noise $d\mathbf{B}$ of variance $(1 \times 10^{-12})^2 \text{ A}^2$ makes the model stochastic, and
 1664 introduces variability in frequency at each parameterization \mathbf{z} .

1665 In order to measure the frequency of the hub neuron during EPI, the STG model was simulated for
 1666 $T = 300$ time steps of $dt = 25\text{ms}$. The chosen dt and T were the most computationally convenient
 1667 choices yielding accurate frequency measurement. We used a basis of complex exponentials with
 1668 frequencies from 0.0-1.0 Hz at 0.01Hz resolution to measure frequency from simulated time series

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

1669 To measure spiking frequency, we processed simulated membrane potentials with a relu (spike
 1670 extraction) and low-pass filter with averaging window of size 20, then took the frequency with the
 1671 maximum absolute value of the complex exponential basis coefficients of the processed time-series.
 1672 The first 20 temporal samples of the simulation are ignored to account for initial transients.

1673 To differentiate through the maximum frequency identification, we used a soft-argmax Let $X_\alpha \in$
 1674 $\mathcal{C}^{|\Phi|}$ be the complex exponential filter bank dot products with the signal $x_\alpha \in \mathbb{R}^N$, where $\alpha \in$
 1675 $\{\text{f1}, \text{f2}, \text{hub}, \text{s1}, \text{s2}\}$. The soft-argmax is then calculated using temperature parameter $\beta = 100$

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

1676 where $i = [0, 1, \dots, 100]$. The frequency is then calculated as

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

1677 Intermediate hub frequency, like all other emergent properties in this work, is defined by the mean
 1678 and variance of the emergent property statistics. In this case, we have one statistic, hub neuron
 1679 frequency, where the mean was chosen to be 0.55Hz, and variance was chosen to be $(0.025\text{Hz})^2$ to
 1680 capture variation in frequency between 0.5Hz and 0.6Hz (Equation 2A deep generative modeling
 1681 approach to emergent property inferenceequation.3.2). As a maximum entropy distribution, $T(\mathbf{x}, \mathbf{z})$
 1682 is comprised of both these first and second moments of the hub neuron frequency (as in Equations
 1683 ?? and ??)

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

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

1684
 1685 Throughout optimization, the augmented Lagrangian parameters η and c , were updated after each
 1686 epoch of 5,000 iterations(see Section 5.1.3Augmented Lagrangian optimizationsubsubsection.5.1.3).
 1687 The optimization converged after five epochs (Fig. S4).

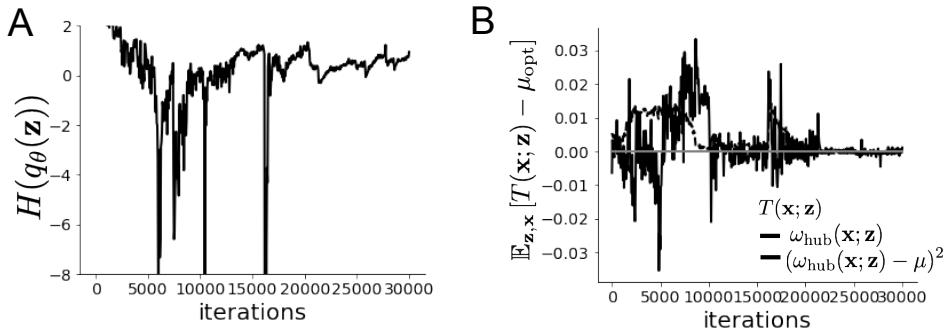


Figure 8: (STG1): EPI optimization of the STG model producing network syncing. A. Entropy throughout optimization. B. The emergent property statistic means and variances converge to their constraints at 25,000 iterations following the fifth augmented Lagrangian epoch.

1688 For EPI in Fig 1 Emergent property inference (EPI) in the stomatogastric ganglion. **A.** Conductance-based
 1689 biophysical model of the STG subcircuit. In the STG model, jagged connections indicate electrical coupling
 1690 having electrical conductance g_{el} . Other connections in the diagram are inhibitory synaptic projections
 1691 having strength g_{synA} onto the hub neuron, and $g_{\text{synB}} = 5\text{nS}$ for mutual inhibitory connections. Parameters
 1692 are represented by the vector \mathbf{z} and membrane potentials by the vector \mathbf{x} . The evolution of this model's
 1693 activity $\mathbf{x}(t)$ is predicated by differential equations. **B.** Spiking frequency $\omega(\mathbf{x}; \mathbf{z})$ is an emergent property
 1694 statistic. In this example, spiking frequency is measured from simulated activity of the STG model at
 1695 parameter choices of $g_{\text{el}} = 4.5\text{nS}$ and $g_{\text{synA}} = 3\text{nS}$. **C.** The emergent property of intermediate hub frequency,
 1696 in which the hub neuron fires at a rate between the fast and slow frequencies. This emergent property is
 1697 defined by a mean and variance on the emergent property statistic. Simulated activity traces are colored
 1698 by log probability density of their generating parameters in the EPI-inferred distribution (Panel E). **D.**
 1699 For a choice of model and emergent property, emergent property inference (EPI) learns a deep probability
 1700 distribution of parameters \mathbf{z} . Deep probability distributions map a simple random variable \mathbf{z}_0 through a
 1701 deep neural network with weights and biases $\boldsymbol{\theta}$ to parameters $\mathbf{z} = g_{\boldsymbol{\theta}}(\mathbf{z}_0)$. In EPI optimization, stochastic
 1702 gradient steps in $\boldsymbol{\theta}$ are taken such that entropy is maximized, and the emergent property \mathcal{X} is produced. The
 1703 EPI posterior distribution is denoted $q_{\boldsymbol{\theta}}(\mathbf{z} | \mathcal{X})$. **E.** The EPI posterior producing intermediate hub frequency.
 1704 Samples are colored by log probability density. Distribution contours of average hub neuron frequency from
 1705 mean of .55 Hz are shown at levels of .525, .53,575 Hz (dark to light gray away from mean). Eigenvectors
 1706 of the Hessian at the mode of the inferred distribution are indicated as \mathbf{v}_1 (solid) and \mathbf{v}_2 (dashed) with lengths
 1707 scaled by the square root of the absolute value of their eigenvalues. **F** Simulations from parameters in E.
 1708 (Top) The predictive distribution of the posterior obeys the emergent property. The black and gray dashed

1709 lines show the mean and two standard deviations according the emergent property, respectively. (Bottom)
1710 Simulations at the starred parameter values figure.1E, we used a real NVP architecture with three Real
1711 NVP coupling layers and two-layer neural networks of 25 units per layer. The normalizing flow
1712 architecture mapped $z_0 \sim \mathcal{N}(\mathbf{0}, I)$ to a support of $\mathbf{z} = [g_{\text{el}}, g_{\text{synA}}] \in [4, 8] \times [0.01, 4]$, initialized to
1713 a gaussian approximation of samples returned by a preliminary ABC search. We did not include
1714 $g_{\text{synA}} < 0.01$, for numerical stability. EPI optimization was run using 5 different random seeds for
1715 architecture initialization $\boldsymbol{\theta}$ with an augmented Lagrangian coefficient of $c_0 = 10^5$, a batch size
1716 $n = 400$, and $\beta = 2$. The distribution shown is that of the architecture converging with criteria
1717 $N_{\text{test}} = 100$ at greatest entropy across random seeds.

1718 We calculated the Hessian at the mode of the inferred EPI distribution. The Hessian of a probability
1719 model is the second order gradient of the log probability density $\log q_{\boldsymbol{\theta}}(\mathbf{z})$ with respect to the
1720 parameters \mathbf{z} : $\frac{\partial^2 \log q_{\boldsymbol{\theta}}(\mathbf{z})}{\partial \mathbf{z} \partial \mathbf{z}^\top}$. With EPI, we can examine the Hessian, which is analytically available
1721 throughout distribution, to indicate the dimensions of parameter space that are sensitive (strongly
1722 negative eigenvalue), and which are degenerate (low magnitude eigenvalue) with respect to the
1723 emergent property produced. In Figure 1Emergent property inference (EPI) in the stomatogastric
1724 ganglion. **A.** Conductance-based biophysical model of the STG subcircuit. In the STG model, jagged
1725 connections indicate electrical coupling having electrical conductance g_{el} . Other connections in the diagram
1726 are inhibitory synaptic projections having strength g_{synA} onto the hub neuron, and $g_{\text{synB}} = 5\text{nS}$ for mutual
1727 inhibitory connections. Parameters are represented by the vector \mathbf{z} and membrane potentials by the vector
1728 \mathbf{x} . The evolution of this model's activity $\mathbf{x}(t)$ is predicated by differential equations. **B.** Spiking frequency
1729 $\omega(\mathbf{x}; \mathbf{z})$ is an emergent property statistic. In this example, spiking frequency is measured from simulated
1730 activity of the STG model at parameter choices of $g_{\text{el}} = 4.5\text{nS}$ and $g_{\text{synA}} = 3\text{nS}$. **C.** The emergent property
1731 of intermediate hub frequency, in which the hub neuron fires at a rate between the fast and slow frequencies.
1732 This emergent property is defined by a mean and variance on the emergent property statistic. Simulated
1733 activity traces are colored by log probability density of their generating parameters in the EPI-inferred
1734 distribution (Panel E). **D.** For a choice of model and emergent property, emergent property inference (EPI)
1735 learns a deep probability distribution of parameters \mathbf{z} . Deep probability distributions map a simple random
1736 variable \mathbf{z}_0 through a deep neural network with weights and biases $\boldsymbol{\theta}$ to parameters $\mathbf{z} = g_{\boldsymbol{\theta}}(\mathbf{z}_0)$. In EPI
1737 optimization, stochastic gradient steps in $\boldsymbol{\theta}$ are taken such that entropy is maximized, and the emergent
1738 property \mathcal{X} is produced. The EPI posterior distribution is denoted $q_{\boldsymbol{\theta}}(\mathbf{z} \mid \mathcal{X})$. **E.** The EPI posterior
1739 producing intermediate hub frequency. Samples are colored by log probability density. Distribution contours
1740 of average hub neuron frequency from mean of .55 Hz are shown at levels of .525, .53,575 Hz (dark

1741 to light gray away from mean). Eigenvectors of the Hessian at the mode of the inferred distribution are
 1742 indicated as \mathbf{v}_1 (solid) and \mathbf{v}_2 (dashed) with lengths scaled by the square root of the absolute value of their
 1743 eigenvalues. **F** Simulations from parameters in E. (Top) The predictive distribution of the posterior obeys the
 1744 emergent property. The black and gray dashed lines show the mean and two standard deviations according
 1745 to the emergent property, respectively. (Bottom) Simulations at the starred parameter values figure.1D, the
 1746 eigenvectors of the Hessian v_1 (solid) and v_2 (dashed) are shown evaluated at the mode of the
 1747 distribution. The length of the arrows is inversely proportional to the square root of absolute
 1748 value of their eigenvalues $\lambda_1 = -10.7$ and $\lambda_2 = -3.22$. Since the Hessian eigenvectors have sign
 1749 degeneracy, the visualized directions in 2-D parameter space are chosen arbitrarily.

1750 5.2.2 Primary visual cortex

1751 In the stochastic stabilized supralinear network [71], population rate responses \mathbf{x} to input \mathbf{h} , recur-
 1752 rent input $W\mathbf{x}$ and slow noise ϵ are governed by

$$\tau \frac{d\mathbf{x}}{dt} = -\mathbf{x} + \phi(W\mathbf{x} + \mathbf{h} + \epsilon), \quad (62)$$

1753 where the noise is an Ornstein-Uhlenbeck process $\epsilon \sim OU(\tau_{\text{noise}}, \sigma)$

$$\tau_{\text{noise}} d\epsilon_\alpha = -\epsilon_\alpha dt + \sqrt{2\tau_{\text{noise}}} \tilde{\sigma}_\alpha dB \quad (63)$$

1754 with $\tau_{\text{noise}} = 5\text{ms} > \tau = 1\text{ms}$. The noisy process is parameterized as

$$\tilde{\sigma}_\alpha = \sigma_\alpha \sqrt{1 + \frac{\tau}{\tau_{\text{noise}}}}, \quad (64)$$

1755 so that σ parameterizes the variance of the noisy input in the absence of recurrent connectivity
 1756 ($W = \mathbf{0}$). As contrast increases, input to the E- and P-populations increases relative to a baseline
 1757 input $\mathbf{h} = \mathbf{h}_b + c\mathbf{h}_c$. Connectivity (W_{fit}) and input ($\mathbf{h}_{b,\text{fit}}$ and $\mathbf{h}_{c,\text{fit}}$) parameters were fit using the
 1758 deterministic V1 circuit model [72]

$$W_{\text{fit}} = \begin{bmatrix} W_{EE} & W_{EP} & W_{ES} & W_{EV} \\ W_{PE} & W_{PP} & W_{PS} & W_{PV} \\ W_{SE} & W_{SP} & W_{SS} & W_{SV} \\ W_{VE} & W_{VP} & W_{VS} & W_{VV} \end{bmatrix} = \begin{bmatrix} 2.18 & -1.19 & -.594 & -.229 \\ 1.66 & -.651 & -.680 & -.242 \\ .895 & -5.22 \times 10^{-3} & -1.51 \times 10^{-4} & -.761 \\ 3.34 & -2.31 & -.254 & -2.52 \times 10^{-4} \end{bmatrix}, \quad (65)$$

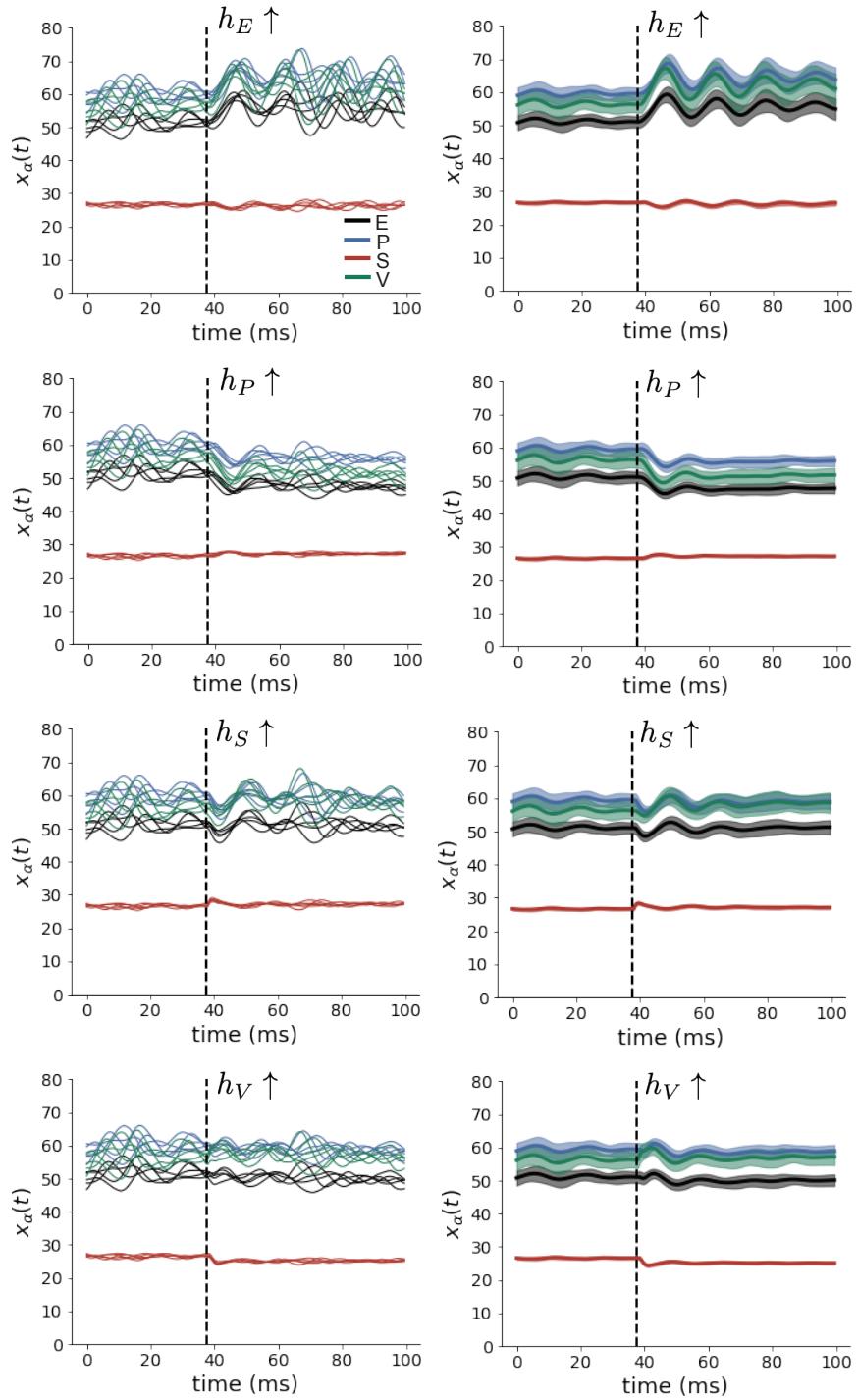


Figure 9: (V1 1) (Left) Simulations for small increases in neuron-type population input. Input magnitudes are chosen so that effect is salient (0.002 for E and P, but 0.02 for S and V). (Right) Average (solid) and standard deviation (shaded) of stochastic fluctuations of responses.

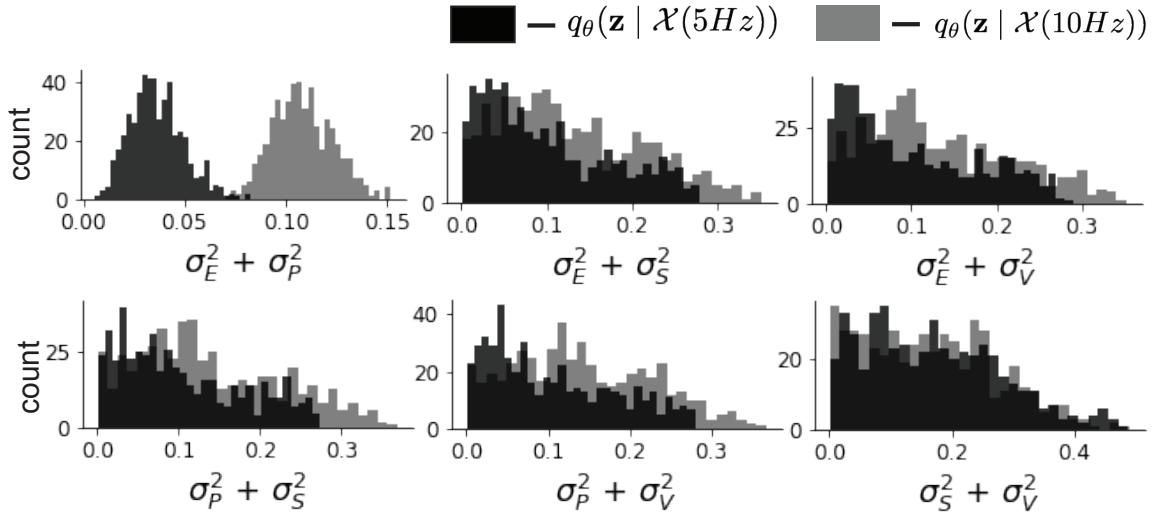


Figure 10: (V1 2) Posterior predictive distributions of the sum of squares of each pair of noise parameters.

$$\mathbf{h}_{b,\text{fit}} = \begin{bmatrix} .416 \\ .429 \\ .491 \\ .486 \end{bmatrix}, \quad (66)$$

1759 and

$$\mathbf{h}_{c,\text{fit}} = \begin{bmatrix} .359 \\ .403 \\ 0 \\ 0 \end{bmatrix}. \quad (67)$$

1760 To obtain rates on a realistic scale (100-fold greater), we map these fitted parameters to an equivalence class
 1761

$$W = \begin{bmatrix} W_{EE} & W_{EP} & W_{ES} & W_{EV} \\ W_{PE} & W_{PP} & W_{PS} & W_{PV} \\ W_{SE} & W_{SP} & W_{SS} & W_{SV} \\ W_{VE} & W_{VP} & W_{VS} & W_{VV} \end{bmatrix} = \begin{bmatrix} .218 & -.119 & -.0594 & -.0229 \\ .166 & -.0651 & -.068 & -.0242 \\ .0895 & -5.22 \times 10^{-4} & -1.51 \times 10^{-5} & -.0761 \\ .334 & -.231 & -.0254 & -2.52 \times 10^{-5} \end{bmatrix}, \quad (68)$$

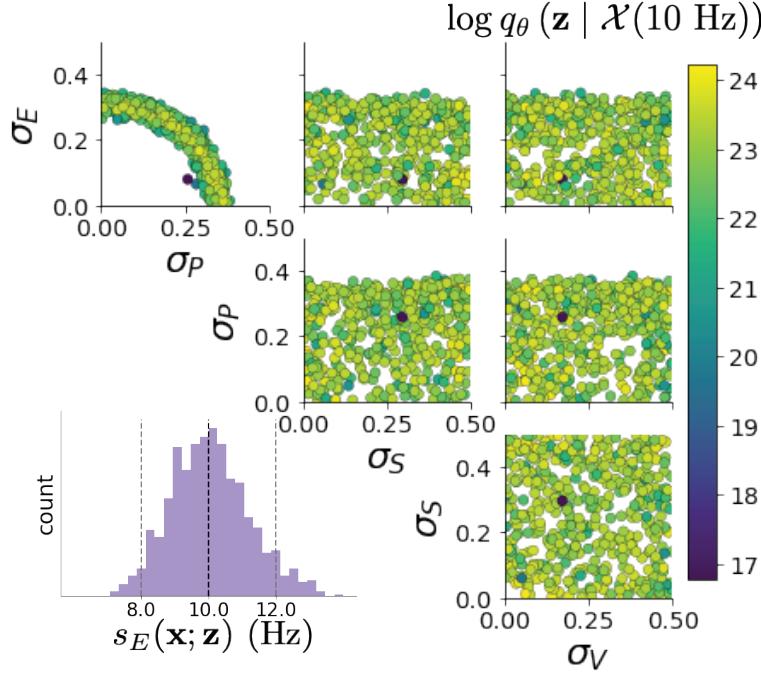


Figure 11: (V1 3) EPI posterior for $\mathcal{X}(10 \text{ Hz})$.

$$\mathbf{h}_b = \begin{bmatrix} h_{b,E} \\ h_{b,P} \\ h_{b,S} \\ h_{b,V} \end{bmatrix} = \begin{bmatrix} 4.16 \\ 4.29 \\ 4.91 \\ 4.86 \end{bmatrix}, \quad (69)$$

1762 and

$$\mathbf{h}_c = \begin{bmatrix} h_{c,E} \\ h_{c,P} \\ h_{c,S} \\ h_{c,V} \end{bmatrix} = \begin{bmatrix} 3.59 \\ 4.03 \\ 0 \\ 0 \end{bmatrix}. \quad (70)$$

1763 Circuit responses are simulated using $T = 200$ time steps at $dt = 0.5\text{ms}$ from an initial condition
 1764 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
 1765 as the square root of the temporal variance from $t_{ss} = 75\text{ms}$ to $Tdt = 100\text{ms}$ averaged over 100
 1766 independent trials.

$$s_E(\mathbf{x}; \mathbf{z}) = \mathbb{E}_x \left[\sqrt{\mathbb{E}_{t > t_{ss}} [(x_E(t) - \mathbb{E}_{t > t_{ss}} [x_E(t)])^2]} \right] \quad (71)$$

1767 For EPI in Fig 2 Emergent property inference in the stochastic stabilized supralinear network (SSSN) A.
 1768 Four-population model of primary visual cortex with excitatory (black), parvalbumin (blue), somatostatin

1769 (red), and VIP (green) neurons (excitatory and inhibitory projections filled and unfilled, respectively). Some
 1770 neuron-types largely do not form synaptic projections to others ($|W_{\alpha_1, \alpha_2}| < 0.025$). Each neural population
 1771 receives a baseline input \mathbf{h}_b , and the E- and P-populations also receive a contrast-dependent input \mathbf{h}_c .
 1772 Additionally, each neural population receives a slow noisy input ϵ . **B.** Steady-state responses of the SSSN
 1773 model (deterministic, $\sigma = \mathbf{0}$) to varying contrasts. The response at 50% contrast (dots) is the focus of our
 1774 analysis. **C.** Transient network responses of the SSSN model at 50 % contrast. (Left) Traces are independent
 1775 trials with varying initialization $\mathbf{x}(0)$ and noise realization. (Right) Mean (solid line) and standard deviation
 1776 (shading) of responses. **D.** EPI posterior of noise parameters \mathbf{z} conditioned on E-population variability. The
 1777 posterior predictive distribution of $s_E(\mathbf{x}; \mathbf{z})$ is show on the bottom-left. **E.** (Top) Enlarged visualization of the
 1778 σ_E - σ_P marginal distribution of the posteriors $q_{\theta}(\mathbf{z} | \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} | \mathcal{X}(10 \text{ Hz})$. Each black dot shows the
 1779 mode at each σ_P . The arrows show the most sensitive dimensions of the Hessian evaluated at these modes. **F.**
 1780 The predictive distributions of $\sigma_E^2 + \sigma_P^2$ of each posterior $q_{\theta}(\mathbf{z} | \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} | \mathcal{X}(10 \text{ Hz}))$ figure.2D-E, we
 1781 used a real NVP architecture with three Real NVP coupling layers and two-layer neural networks
 1782 of 50 units per layer. The normalizing flow architecture mapped $z_0 \sim \mathcal{N}(\mathbf{0}, I)$ to a support of
 1783 $\mathbf{z} = [\sigma_E, \sigma_P, \sigma_S, \sigma_V] \in [0.0, 0.5]^4$. EPI optimization was run using three different random seeds for
 1784 architecture initialization θ with an augmented Lagrangian coefficient of $c_0 = 10^{-1}$, a batch size
 1785 $n = 100$, and $\beta = 2$. The distributions shown are those of the architectures converging with criteria
 1786 $N_{\text{test}} = 100$ at greatest entropy across random seeds.

1787 In Fig. 2Emergent property inference in the stochastic stabilized supralinear network (SSSN) **A.** Four-
 1788 population model of primary visual cortex with excitatory (black), parvalbumin (blue), somatostatin (red),
 1789 and VIP (green) neurons (excitatory and inhibitory projections filled and unfilled, respectively). Some
 1790 neuron-types largely do not form synaptic projections to others ($|W_{\alpha_1, \alpha_2}| < 0.025$). Each neural population
 1791 receives a baseline input \mathbf{h}_b , and the E- and P-populations also receive a contrast-dependent input \mathbf{h}_c .
 1792 Additionally, each neural population receives a slow noisy input ϵ . **B.** Steady-state responses of the SSSN
 1793 model (deterministic, $\sigma = \mathbf{0}$) to varying contrasts. The response at 50% contrast (dots) is the focus of our
 1794 analysis. **C.** Transient network responses of the SSSN model at 50 % contrast. (Left) Traces are independent
 1795 trials with varying initialization $\mathbf{x}(0)$ and noise realization. (Right) Mean (solid line) and standard deviation
 1796 (shading) of responses. **D.** EPI posterior of noise parameters \mathbf{z} conditioned on E-population variability. The
 1797 posterior predictive distribution of $s_E(\mathbf{x}; \mathbf{z})$ is show on the bottom-left. **E.** (Top) Enlarged visualization of the
 1798 σ_E - σ_P marginal distribution of the posteriors $q_{\theta}(\mathbf{z} | \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} | \mathcal{X}(10 \text{ Hz})$. Each black dot shows the
 1799 mode at each σ_P . The arrows show the most sensitive dimensions of the Hessian evaluated at these modes.
 1800 **F.** The predictive distributions of $\sigma_E^2 + \sigma_P^2$ of each posterior $q_{\theta}(\mathbf{z} | \mathcal{X}(5 \text{ Hz})$ and $q_{\theta}(\mathbf{z} | \mathcal{X}(10 \text{ Hz}))$ figure.2E,

1801 we visualize the modes of $q_{\theta}(\mathbf{z} \mid \mathcal{X})$ throughout the σ_E - σ_P marginal. Specifically, we calculated

$$\begin{aligned}\mathbf{z}^*(\sigma_{P,\text{fixed}}) &= \underset{\mathbf{z}}{\operatorname{argmax}} \log q_{\theta}(\mathbf{z} \mid \mathcal{X}) \\ \text{s.t. } \sigma_P &= \sigma_{P,\text{fixed}}\end{aligned}\tag{72}$$

1802 At each mode \mathbf{z}^* , we calculated the Hessian and visualized the sensitivity dimension in the direction
1803 of positive σ_E .

1804 **5.2.3 Primary visual cortex: challenges to analysis**

1805 TODO Agostina and I are putting this together now.

1806 **5.2.4 Superior colliculus**

1807 In the model of Duan et al [55], there are four total units: two in each hemisphere corresponding to
1808 the Pro/Contra and Anti/Ipsi populations. They are denoted as left Pro (LP), left Anti (LA), right
1809 Pro (RP) and right Anti (RA). Each unit has an activity (x_α) and internal variable (u_α) related
1810 by

$$x_\alpha = \phi(u_\alpha) = \left(\frac{1}{2} \tanh \left(\frac{u_\alpha - a}{b} \right) + \frac{1}{2} \right) \tag{73}$$

1811 where $\alpha \in \{LP, LA, RA, RP\}$, $a = 0.05$ and $b = 0.5$ control the position and shape of the nonlin-
1812 earity, respectively. During periods of optogenetic inactivation, activity was decreased proportional
1813 to the optogenetic strength γ

$$x_\alpha = (1 - \gamma)\phi(u_\alpha). \tag{74}$$

1814 We order the neural populations of x and u in the following manner

$$\mathbf{x} = \begin{bmatrix} x_{LP} \\ x_{LA} \\ x_{RP} \\ x_{RA} \end{bmatrix} \quad \mathbf{u} = \begin{bmatrix} u_{LP} \\ u_{LA} \\ u_{RP} \\ u_{RA} \end{bmatrix}, \tag{75}$$

1815 which evolve according to

$$\tau \frac{d\mathbf{u}}{dt} = -\mathbf{u} + W\mathbf{x} + \mathbf{h} + d\mathbf{B}. \tag{76}$$

1816 with time constant $\tau = 0.09s$, step size 24ms and Gaussian noise $d\mathbf{B}$ of variance 0.2^2 . The weight

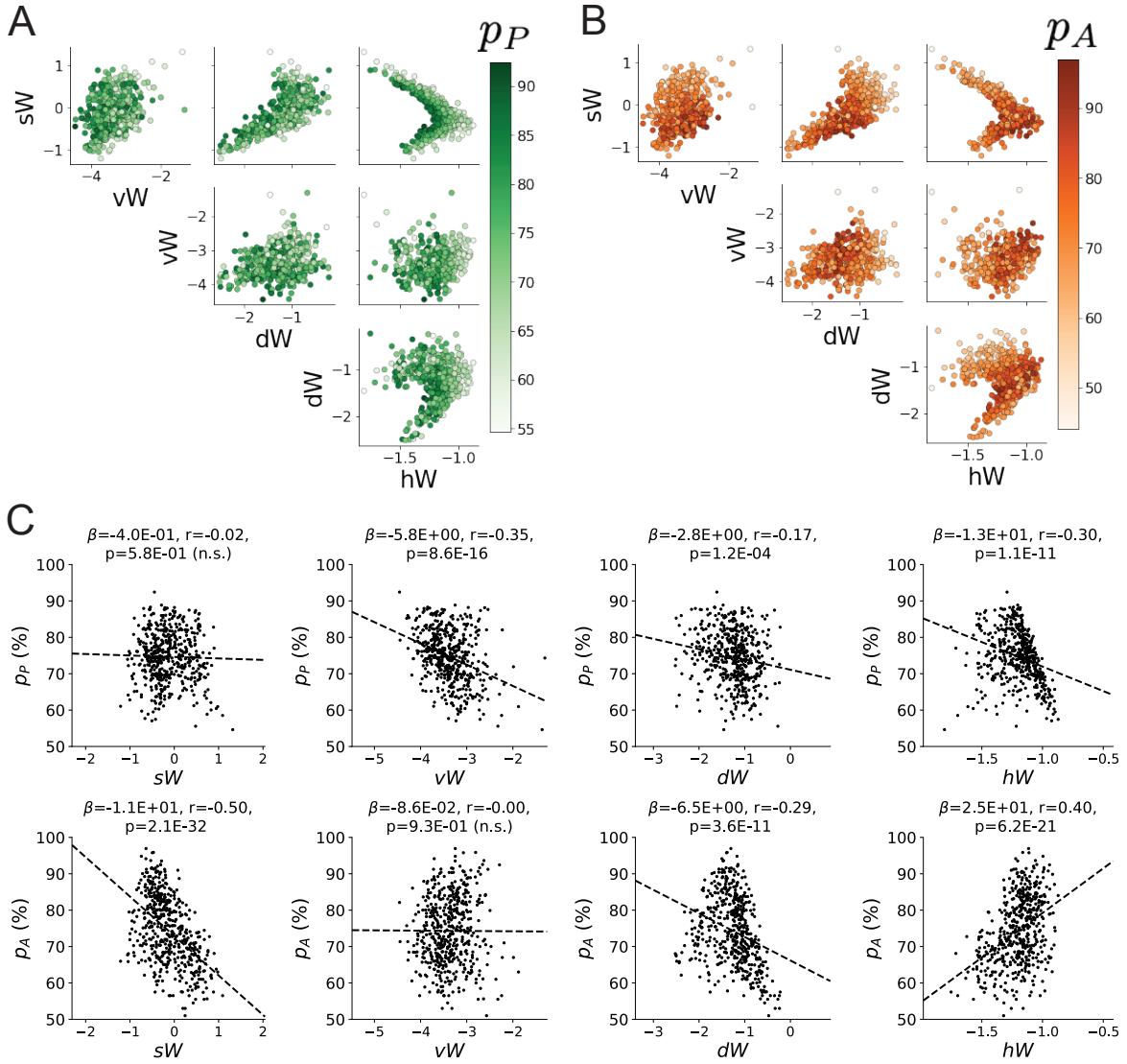


Figure 12: (SC1): **A.** Same pairplot as Fig. 3A. Rapid task switching behavioral paradigm (see text). **B.** Model of superior colliculus (SC). Neurons: LP - left pro, RP - right pro, LA - left anti, RA - right anti. Parameters: sW - self, hW - horizontal, vW - vertical, dW - diagonal weights. **C.** The EPI posterior distribution of rapid task switching networks. Red and purple stars (\mathbf{z}_1 and \mathbf{z}_2) indicate different connectivity regimes with different sensitivity vectors \mathbf{v}_1 and \mathbf{v}_2 . (Middle-left) Posterior predictive distribution of task accuracies. (Bottom-left) Task accuracy along dimensions of sensitivity in each connectivity regime. **D.** Means (solid) and standard deviations (shaded) of each population across random simulated trials. Top plots show Pro (top) and Anti (bottom) responses for connectivity \mathbf{z}_1 . Bottom rows show the same \mathbf{z}_2 . **E.** The EPI posterior predicts experimental results (left) showing no change in the Pro task, but larger error in the Anti task (right). **F.** Accuracy in the Anti task during delay period optogenetic inactivation $p_{A,\text{opto}}$ is strongly anticorrelated with accuracy in the Pro task. **G.** Accuracy with delay period inactivation along each connectivity regime's dimension of sensitivity figure.3C colored by Pro task accuracy. **B.** Same as A colored by Anti task accuracy. **C.** Connectivity parameters of EPI distributions versus task accuracies. β is slope coefficient of linear regression, r is correlation, and p is the two-tailed p-value.

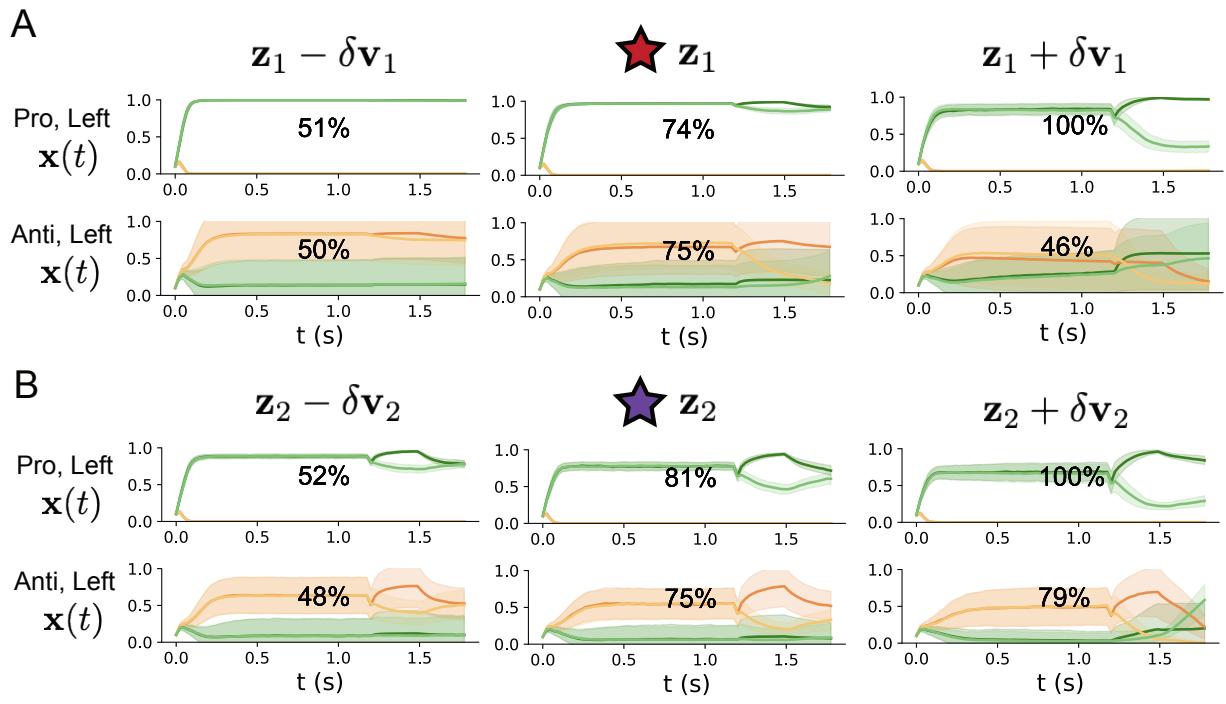


Figure 13: (SC2): **A.** Simulations in network regime \mathbf{z}_1 (center) with simulations given connectivity perturbations in the negative direction of the sensitivity vector \mathbf{v}_1 (left) and positive direction (right). **B.** Same as A for network regime \mathbf{z}_2 .

1817 matrix has 4 parameters sW , vW , hW , and dW :

$$W = \begin{bmatrix} sW & vW & hW & dW \\ vW & sW & dW & hW \\ hW & dW & sW & vW \\ dW & hW & vW & sW \end{bmatrix}. \quad (77)$$

1818 The circuit receives four different inputs throughout each trial, which has a total length of 1.8s.

$$\mathbf{h} = \mathbf{h}_{\text{constant}} + \mathbf{h}_{\text{P,bias}} + \mathbf{h}_{\text{rule}} + \mathbf{h}_{\text{choice-period}} + \mathbf{h}_{\text{light}}. \quad (78)$$

1819 There is a constant input to every population,

$$\mathbf{h}_{\text{constant}} = I_{\text{constant}}[1, 1, 1, 1]^\top, \quad (79)$$

1820 a bias to the Pro populations

$$\mathbf{h}_{\text{P,bias}} = I_{\text{P,bias}}[1, 0, 1, 0]^\top, \quad (80)$$

1821 rule-based input depending on the condition

$$\mathbf{h}_{\text{P,rule}}(t) = \begin{cases} I_{\text{P,rule}}[1, 0, 1, 0]^\top, & \text{if } t \leq 1.2s \\ 0, & \text{otherwise} \end{cases} \quad (81)$$

1822

$$\mathbf{h}_{\text{A,rule}}(t) = \begin{cases} I_{\text{A,rule}}[0, 1, 0, 1]^\top, & \text{if } t \leq 1.2s \\ 0, & \text{otherwise} \end{cases}, \quad (82)$$

1823 a choice-period input

$$\mathbf{h}_{\text{choice}}(t) = \begin{cases} I_{\text{choice}}[1, 1, 1, 1]^\top, & \text{if } t > 1.2s \\ 0, & \text{otherwise} \end{cases}, \quad (83)$$

1824 and an input to the right or left-side depending on where the light stimulus is delivered

$$\mathbf{h}_{\text{light}}(t) = \begin{cases} I_{\text{light}}[1, 1, 0, 0]^\top, & \text{if } 1.2s < t < 1.5s \text{ and Left} \\ I_{\text{light}}[0, 0, 1, 1]^\top, & \text{if } 1.2s < t < 1.5s \text{ and Right} \\ 0, & \text{otherwise} \end{cases}. \quad (84)$$

1825 The input parameterization was fixed to $I_{\text{constant}} = 0.75$, $I_{\text{P,bias}} = 0.5$, $I_{\text{P,rule}} = 0.6$, $I_{\text{A,rule}} = 0.6$,

1826 $I_{\text{choice}} = 0.25$, and $I_{\text{light}} = 0.5$.

1827 The accuracies of p_P and p_A are calculated as

$$p_P(\mathbf{x}; \mathbf{z}) = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}|\mathbf{z})} [\Theta[x_{LP}(t = 1.8s) - x_{RP}(t = 1.8s)]] \quad (85)$$

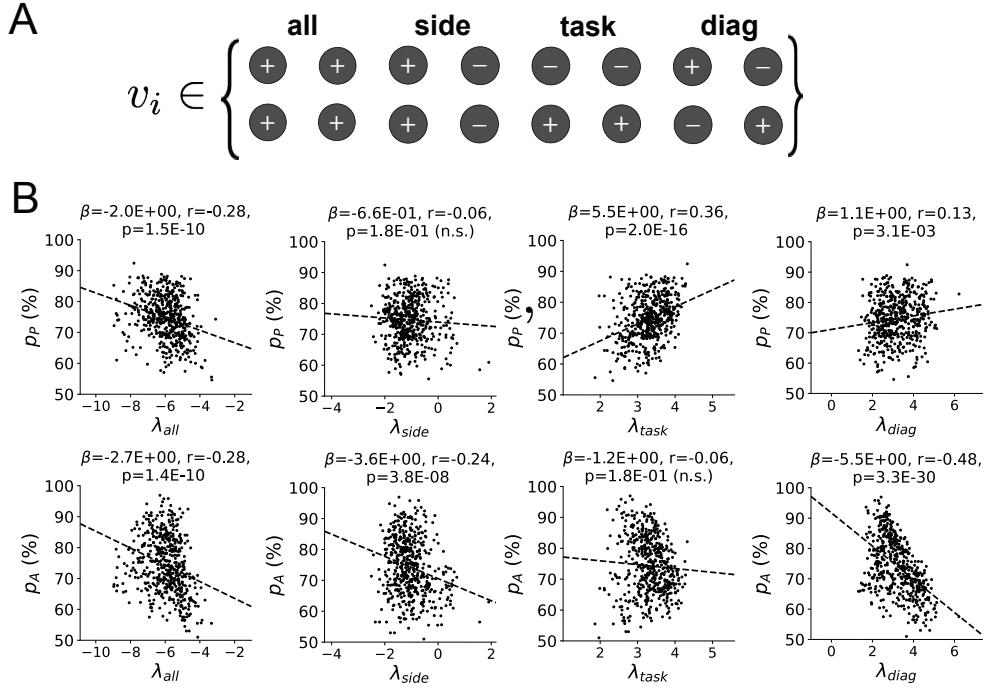


Figure 14: (SC3): **A.** Invariant eigenvectors of connectivity matrix W . **B.** Eigenvalues of connectivities of EPI distribution versus task accuracies.

1828 and

$$p_A(\mathbf{x}; \mathbf{z}) = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}|\mathbf{z})} [\Theta[x_{RP}(t = 1.8s) - x_{LP}(t = 1.8s)]] \quad (86)$$

1829 given that the stimulus is on the left side, where Θ is the Heaviside step function, and the accuracy
1830 is averaged over 200 independent trials. The Heaviside step function is approximated as

$$\Theta(\mathbf{x}) = \text{sigmoid}(\beta\mathbf{x}), \quad (87)$$

1831 where $\beta = 100$.

1832 Writing the EPI posterior as a maximum entropy distribution, $T(\mathbf{x}, \mathbf{z})$ is comprised of both these
1833 first and second moments of the accuracy in each task (as in Equations ?? and ??)

$$T(\mathbf{x}; \mathbf{z}) = \begin{bmatrix} p(\mathbf{x}; \mathbf{z})_P \\ p(\mathbf{x}; \mathbf{z})_A \\ (p(\mathbf{x}; \mathbf{z})_P - 75\%)^2 \\ (p(\mathbf{x}; \mathbf{z})_A - 75\%)^2 \end{bmatrix}, \quad (88)$$

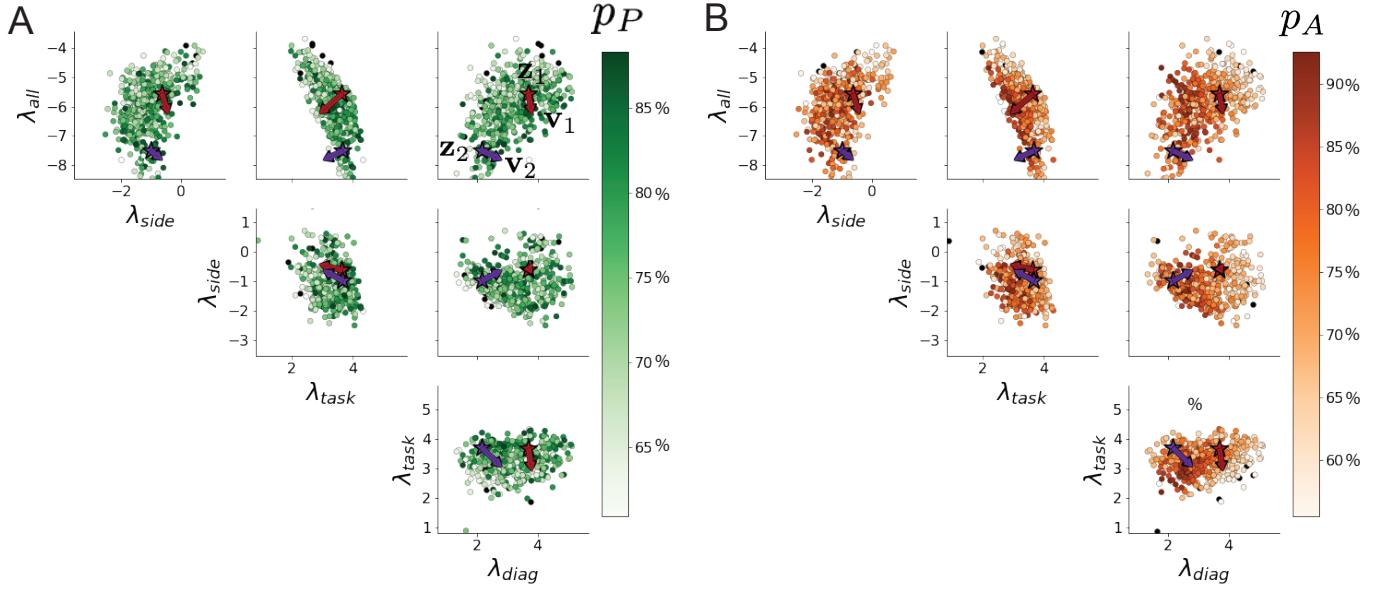


Figure 15: (SC4): **A.** Pairplots of eigenvalues of connectivity matrices in EPI distribution colored by Pro task accuracy. Red and purple stars and arrows correspond to eigenvalues and sensitivity directions \mathbf{z}_1 , \mathbf{z}_2 , \mathbf{v}_1 , and \mathbf{v}_2 . **B.** Same colored by Anti task accuracy.

1834

$$\boldsymbol{\mu}_{\text{opt}} = \begin{bmatrix} 75\% \\ 75\% \\ 7.5\%^2 \\ 7.5\%^2 \end{bmatrix}. \quad (89)$$

1835 Throughout optimization, the augmented Lagrangian parameters η and c , were updated after each
 1836 epoch of 2,000 iterations(see Section 5.1.3Augmented Lagrangian optimizationsubsubsection.5.1.3).

1837 The optimization converged after six epochs (Fig. 17(SC6): A. Entropy throughout optimization. B.

1838 The emergent property statistic means and variances converge to their constraints at 20,000 iterations
 1839 following the tenth augmented Lagrangian epochfigure.17).

1840 For EPI in Fig. 3**A**. Rapid task switching behavioral paradigm (see text). **B.** Model of superior colliculus (SC).

1841 Neurons: LP - left pro, RP - right pro, LA - left anti, RA - right anti. Parameters: sW - self, hW - horizontal, vW

1842 -vertical, dW - diagonal weights. **C.** The EPI posterior distribution of rapid task switching networks. Red and purple

1843 stars (\mathbf{z}_1 and \mathbf{z}_2) indicate different connectivity regimes with different sensitivity vectors \mathbf{v}_1 and \mathbf{v}_2 . (Middle-left)

1844 Posterior predictive distribution of task accuracies. (Bottom-left) Task accuracy along dimensions of sensitivity in

1845 each connectivity regime. **D.** Means (solid) and standard deviations (shaded) of each population across random

1846 simulated trials. Top plots show Pro (top) and Anti (bottom) responses for connectivity \mathbf{z}_1 . Bottom rows show the

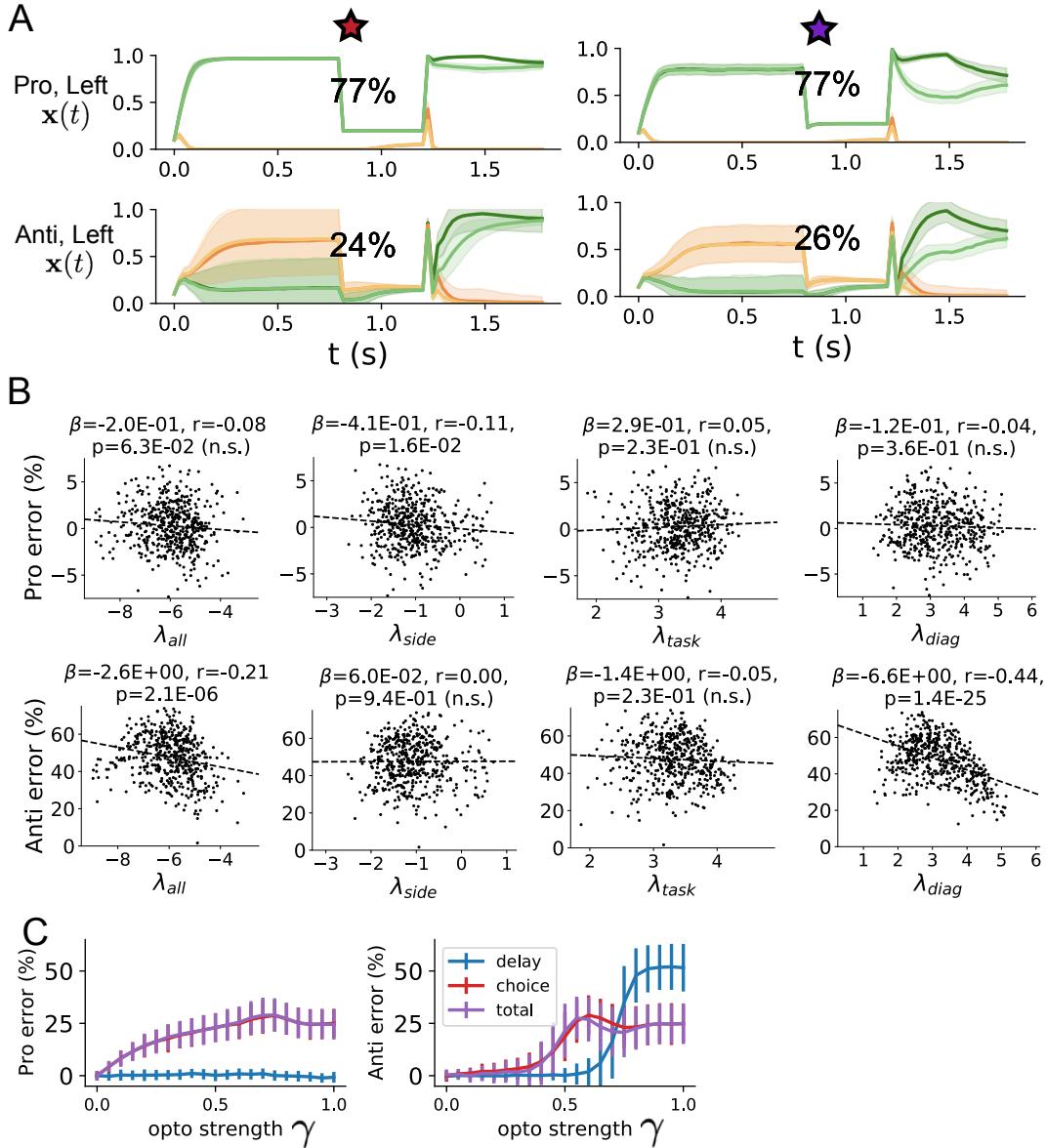


Figure 16: (SC5): **A.** Response of each parameter regime to optogenetic silencing during the delay period. **B.** Connectivity eigenvalues versus the task error induced by delay period inactivation. **C.** Error induced by delay period inactivation with increasing optogenetic strength. Means and standard deviations are calculated across the entire EPI posterior.

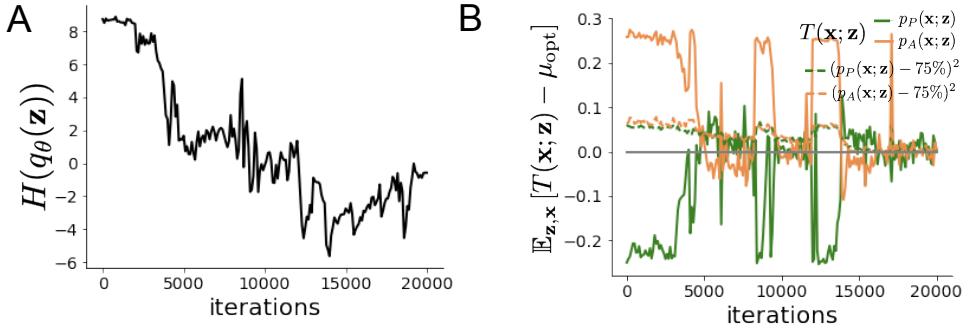


Figure 17: (SC6): A. Entropy throughout optimization. B. The emergent property statistic means and variances converge to their constraints at 20,000 iterations following the tenth augmented Lagrangian epoch.

1847 same \mathbf{z}_2 . **E**. The EPI posterior predicts experimental results (left) showing no change in the Pro task, but larger error
 1848 in the Anti task (right). **F**. Accuracy in the Anti task during delay period optogenetic inactivation $p_{A,\text{opto}}$ is strongly
 1849 anticorrelated with accuracy in the Pro task. **G**. Accuracy with delay period inactivation along each connectivity
 1850 regime's dimension of sensitivityfigure.3C, we used a real NVP architecture with three coupling layers of
 1851 affine transformations parameterized by two-layer neural networks of 50 units per layer. The initial
 1852 distribution was a standard isotropic gaussian $z_0 \sim \mathcal{N}(\mathbf{0}, I)$ mapped to a support of $\mathbf{z}_i \in [-5, 5]$.
 1853 We used an augmented Lagrangian coefficient of $c_0 = 10^2$, a batch size $n = 100$, and $\beta = 2$. The
 1854 distribution shown is that of the architecture converging with criteria $N_{\text{test}} = 25$ at greatest entropy
 1855 across random seeds.

1856 To make sense of this inferred distribution, we identified two modes used to represent the two
 1857 regimes of connectivity in this posterior:

$$\begin{aligned} \mathbf{z}_1 &= \underset{\mathbf{z}}{\operatorname{argmax}} \log q_\theta(\mathbf{z} \mid \mathcal{X}) \\ \text{s.t. } hw &= -1.25, sW > 0 \end{aligned} \tag{90}$$

1858 and

$$\begin{aligned} \mathbf{z}_2 &= \underset{\mathbf{z}}{\operatorname{argmax}} \log q_\theta(\mathbf{z} \mid \mathcal{X}) \\ \text{s.t. } hw &= -1.25, sW < 0 \end{aligned} \tag{91}$$

1859 To understand the connectivity mechanisms governing task accuracy, we took the eigendecomposi-
 1860 tion of the symmetric connectivity matrices $W = V\Lambda V^{-1}$, which results in the same basis vectors
 1861 \mathbf{v}_i for all W parameterized by \mathbf{z} (Fig. 14(SC3): **A**. Invariant eigenvectors of connectivity matrix W . **B**.
 1862 Eigenvalues of connectivities of EPI distribution versus task accuraciesfigure.14A). These basis vectors
 1863 have intuitive roles in processing for this task, and are accordingly named the *all* mode - all neurons

1864 co-fluctuate, *side* mode - one side dominates the other, *task* mode - the Pro or Anti populations
1865 dominate the other, and *diag* mode - Pro- and Anti-populations of opposite hemispheres dominate
1866 the opposite pair. We found significant trends across the EPI posterior connectivities: the eigen-
1867 values λ_{task} and λ_{diag} were correlated with p_P , while λ_{all} was anticorrelated with p_P . λ_{all} , λ_{side} ,
1868 and λ_{diag} were all significantly anticorrelated with p_A .

1869 Under this decomposition, we can re-visualize the posterior in eigenvalue space (Fig. 15(SC4): **A**.
1870 Pairplots of eigenvalues of connectivity matrices in EPI distribution colored by Pro task accuracy. Red and
1871 purple stars and arrows correspond to eigenvalues and sensitivity directions \mathbf{z}_1 , \mathbf{z}_2 , \mathbf{v}_1 , and \mathbf{v}_2 . **B**. Same
1872 colored by Anti task accuracyfigure.15). Furthermore, we can project the dimensions of sensitivity into
1873 eigenvalue space as well, giving us a more intuitive sense of how connectivity affects computation
1874 in each regime. We see that sensitivity dimensions \mathbf{v}_1 and \mathbf{v}_2 , which cause p_P to increase and
1875 a regime dependent change in p_A , both point in the direction of increasing λ_{side} and decreasing
1876 λ_{task} . These eigenvalue changes are evident in the simulations of connectivity perturbations away
1877 from the modes (Fig. 13(SC2): **A**. Simulations in network regime \mathbf{z}_1 (center) with simulations given
1878 connectivity perturbations in the negative direction of the sensitivity vector \mathbf{v}_1 (left) and positive direction
1879 (right). **B**. Same as A for network regime \mathbf{z}_2 figure.13). As the component of connectivity along \mathbf{v}_1 and
1880 \mathbf{v}_2 becomes stronger (left-to-right), there is less separation between Pro an Anti populations (lower
1881 λ_{task}) and greater separation between Left and Right populations following stimulus presentation
1882 (greater λ_{side}). A key differentiating factor is that \mathbf{v}_1 substantially increases λ_{diag} , while \mathbf{v}_2 does
1883 not.

1884 During optogenetic silencing simulations, activations $x_\alpha(t)$ were set to a fraction of their values ($1 -$
1885 γ), where γ is the optogenetic perturbation strength. We found that λ_{all} and λ_{diag} were significantly
1886 anticorrelated with Anti error during delay period inactivation. Delay period inactivation was from
1887 $0.8 < t < 1.2$, choice period inactivation was for $t > 1.2$ and total inactivation was for the entire
1888 trial.

1889 5.2.5 Rank-2 RNN

1890 Traditional approaches to likelihood-free inference – approximate Bayesian computation (ABC)
1891 methods – randomly sample parameters \mathbf{z} until a suitable set is obtained. State-of-the-art ABC
1892 methods leverage sequential Monte Carlo (SMC) sampling techniques to obtain parameter sets more
1893 efficiently. To obtain more parameter samples, SMC-ABC must be run from scratch again. ABC
1894 methods do not confer log probabilities of samples. Like EPI, sequential neural posterior estimation

1895 (SNPE) uses deep learning to produce flexible posterior approximations. Like traditional Bayesian
1896 inference methods, SNPE conditions directly on the statistics of data. This differs from EPI, where
1897 posteriors are conditioned on emergent properties (moment constraints on the posterior predictive
1898 distribution). Peculiarities of SNPE (density estimation approach, two deep networks) make scaling
1899 in \mathbf{z} prohibitive.

1900 SMC-ABC has many hyperparameters, of which pyABC selects automatically by running some ini-
1901 tial diagnostics upon initialization. In concurrence with the literature, SMC-ABC fails to converge
1902 around 25-30 dimensions, since it's proposal samples never get close enough to the target statis-
1903 tics. We searched over many SNPE hyperparameter choices: $n_{\text{train}} \in [2,000, 10,000, 100,000]$ is the
1904 number of simulations run per training epoch, and $n_{\text{mades}} \in [2, 3]$ is the number of masked autore-
1905 gressive density estimators in the deep parameter distribution architecture. The greater n_{train} , the
1906 longer each epoch will take, but the more likely SNPE may converge during that epoch. Greater
1907 n_{mades} increases the flexibility of the deep parameter distribution of SNPE, but slows optimization.
1908 For the timing plot, we show the fastest among all of these choices, and for the convergence plot,
1909 we show the best convergence among all of these choices. During optimization, we used $n_{\text{atom}}=100$
1910 atomic proposals as is recommended.