

# Interrogating theoretical models of neural computation with deep inference

Sean R. Bittner, Agostina Palmigiano, Alex T. Piet, Chunyu A. Duan, Carlos D. Brody,  
Kenneth D. Miller, and John P. Cunningham.

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

<sup>2</sup> A cornerstone of theoretical neuroscience is the circuit model: a system of equations that captures a  
<sup>3</sup> hypothesized neural mechanism. Such models are valuable when they give rise to an experimentally  
<sup>4</sup> observed phenomenon – whether behavioral or in terms of neural activity – and thus can offer  
<sup>5</sup> insights into neural computation. The operation of these circuits, like all models, critically depends  
<sup>6</sup> on the choices of model parameters. Historically, the gold standard has been to analytically derive  
<sup>7</sup> the relationship between model parameters and computational properties. However, this enterprise  
<sup>8</sup> quickly becomes infeasible as biologically realistic constraints are included into the model increasing  
<sup>9</sup> its complexity, often resulting in *ad hoc* approaches to understanding the relationship between  
<sup>10</sup> model and computation. We bring recent machine learning techniques – the use of deep generative  
<sup>11</sup> models for probabilistic inference – to bear on this problem, learning distributions of parameters  
<sup>12</sup> that produce the specified properties of computation. Importantly, the techniques we introduce offer  
<sup>13</sup> a principled means to understand the implications of model parameter choices on computational  
<sup>14</sup> properties of interest. We motivate this methodology with a worked example analyzing sensitivity in  
<sup>15</sup> the stomatogastric ganglion. We then use it to generate insights into neuron-type input-responsivity  
<sup>16</sup> in a model of primary visual cortex, a new understanding of rapid task switching in superior  
<sup>17</sup> colliculus models, and attribution of bias in recurrent neural networks solving a toy mathematical  
<sup>18</sup> problem. More generally, this work suggests a departure from realism vs tractability considerations,  
<sup>19</sup> towards the use of modern machine learning for sophisticated interrogation of biologically relevant  
<sup>20</sup> models.

## <sup>21</sup> 2 Introduction

<sup>22</sup> The fundamental practice of theoretical neuroscience is to use a mathematical model to understand  
<sup>23</sup> neural computation, whether that computation enables perception, action, or some intermediate  
<sup>24</sup> processing [1]. A neural computation is systematized with a set of equations – the model – and  
<sup>25</sup> these equations are motivated by biophysics, neurophysiology, and other conceptual considerations.  
<sup>26</sup> The function of this system is governed by the choice of model parameters, which when configured

27 in a particular way, give rise to a measurable signature of a computation. The work of analyzing a  
28 model then requires solving the inverse problem: given a computation of interest, how can we reason  
29 about these particular parameter configurations? The inverse problem is crucial for reasoning about  
30 likely parameter values, uniquenesses and degeneracies, attractor states and phase transitions, and  
31 predictions made by the model.

32 Consider the idealized practice: one carefully designs a model and analytically derives how model  
33 parameters govern the computation. Seminal examples of this gold standard include our field’s  
34 understanding of memory capacity in associative neural networks [2] and chaos and autocorrelation  
35 timescales in random neural networks [3] (which use models and analyses originating in physics),  
36 as well as the paradoxical effect in excitatory/inhibitory networks [4]. Unfortunately, as circuit  
37 models include more biological realism, theory via analytic derivation becomes intractable. This  
38 creates an unfavorable tradeoff. On the one hand, one may tractably analyze systems of equations  
39 with unrealistic assumptions (for example symmetry or gaussianity), producing accurate inferences  
40 about parameters of a too-simple model. On the other hand, one may choose a more biologically  
41 accurate, scientifically relevant model at the cost of *ad hoc* approaches to analysis (such as simply  
42 examining simulated activity), potentially resulting in bad inferences and thus erroneous scientific  
43 predictions or conclusions.

44 Of course, this same tradeoff has been confronted in many scientific fields characterized by the  
45 need to do inference in complex models. In response, the machine learning community has made  
46 remarkable progress in recent years, via the use of deep neural networks as a powerful inference  
47 engine: a flexible function family that can map observed phenomena (in this case the measurable  
48 signal of some computation) back to probability distributions quantifying the likely parameter  
49 configurations. One celebrated example of this approach from machine learning, of which we  
50 draw key inspiration for this work, is the variational autoencoder [5, 6], which uses a deep neural  
51 network to induce an (approximate) posterior distribution on hidden variables in a latent variable  
52 model, given data. Indeed, these tools have been used to great success in neuroscience as well,  
53 in particular for interrogating parameters (sometimes treated as hidden states) in models of both  
54 cortical population activity [7, 8, 9, 10] and animal behavior [11, 12, 13]. These works have used  
55 deep neural networks to expand the expressivity and accuracy of statistical models of neural data  
56 [14].

57 However, these inference tools have not significantly influenced the study of theoretical neuroscience  
58 models, for at least three reasons. First, at a practical level, the nonlinearities and dynamics of

many theoretical models are such that conventional inference tools typically produce a narrow set of insights into these models. Indeed, only in the last few years has deep learning research advanced to a point of relevance to this class of problem. Second, the object of interest from a theoretical model is not typically data itself, but rather a qualitative phenomenon – inspection of model behavior, or better, a measurable signature of some computation – an *emergent property* of the model. Third, because theoreticians work carefully to construct a model that has biological relevance, such a model as a result often does not fit cleanly into the framing of a statistical model. Technically, because many such models stipulate a noisy system of differential equations that can only be sampled or realized through forward simulation, they lack the explicit likelihood and priors central to the probabilistic modeling toolkit.

To address these three challenges, we developed an inference methodology – ‘emergent property inference’ – which learns a distribution over parameter configurations in a theoretical model. This distribution has two critical properties: (*i*) it is chosen such that draws from the distribution (parameter configurations) correspond to systems of equations that give rise to a specified emergent property (a set of constraints); and (*ii*) it is chosen to have maximum entropy given those constraints, such that we identify all likely parameters and can use the distribution to reason about parametric sensitivity and degeneracies [15]. First, we stipulate a bijective deep neural network that induces a flexible family of probability distributions over model parameterizations with a probability density we can calculate [16, 17, 18]. Second, we quantify the notion of emergent properties as a set of moment constraints on datasets generated by the model. Thus, an emergent property is not a single data realization, but a phenomenon or a feature of the model, which is ultimately the object of interest in theoretical neuroscience. Conditioning on an emergent property requires a variant of deep probabilistic inference methods, which we have previously introduced [19]. Third, because we cannot assume the theoretical model has explicit likelihood on data or the emergent property of interest, we use stochastic gradient techniques in the spirit of likelihood free variational inference [20]. Taken together, emergent property inference (EPI) provides a methodology for inferring parameter configurations consistent with a particular emergent phenomena in theoretical models. We use a classic example of parametric degeneracy in a biological system, the stomatogastric ganglion [21], to motivate and clarify the technical details of EPI.

Equipped with this methodology, we then investigated three models of current importance in theoretical neuroscience. These models were chosen to demonstrate generality through ranges of biological realism (from conductance-based biophysics to recurrent neural networks), neural system

function (from pattern generation to abstract cognitive function), and network scale (from four to infinite neurons). First, we use EPI to produce a set of verifiable hypotheses of input-responsivity in a four neuron-type dynamical model of primary visual cortex; we then validate these hypotheses in the model. Second, we demonstrated how the systematic application of EPI to levels of task performance can generate experimentally testable hypotheses regarding connectivity in superior colliculus. Third, we use EPI to uncover the sources of bias in a low-rank recurrent neural network executing a toy mathematical computation. The novel scientific insights offered by EPI contextualize and clarify the previous studies exploring these models [22, 23, 24, 25] and more generally, suggests a departure from realism vs tractability considerations towards the use of modern machine learning for sophisticated interrogation of biologically relevant models.

We note that, during our preparation and early presentation of this work [26, 27], another work has arisen with broadly similar goals: bringing statistical inference to mechanistic models of neural circuits [28]. We are excited by this broad problem being recognized by the community, and we emphasize that these works offer complementary neuroscientific contributions and use different technical methodologies. While we have advanced our research on deep generative modeling [19] to a point of significant relevance to statistical inference in theoretical neuroscience, they have also furthered their research on approximate Bayesian inference in such models [?]. The existence of these complementary methodologies emphasizes the increased importance and timeliness of both works.

## 3 Results

### 3.1 Motivating emergent property inference of theoretical models

Consideration of the typical workflow of theoretical modeling clarifies the need for emergent property inference. First, one designs or chooses an existing model that, it is hypothesized, captures the computation of interest. To ground this process in a well-known example, consider the stomatogastric ganglion (STG) of crustaceans, a small neural circuit which generates multiple rhythmic muscle activation patterns for digestion [29]. Despite full knowledge of STG connectivity and a precise characterization of its rhythmic pattern generation, biophysical models of the STG have complicated relationships between circuit parameters and neural activity [21, 30]. A model of the STG [22] is shown schematically in Figure 1A, and note that the behavior of this model will be critically dependent on its parameterization – the choices of conductance parameters  $z = [g_{el}, g_{synA}]$ .

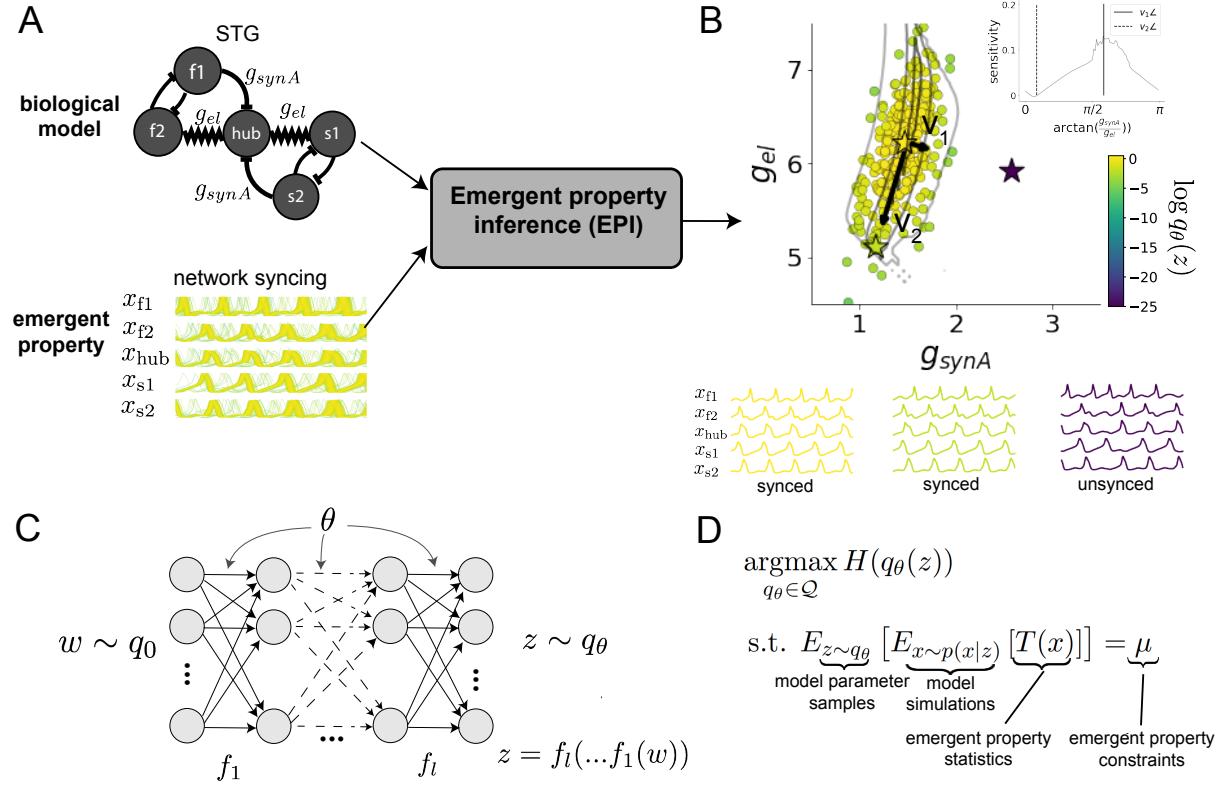


Figure 1: Emergent property inference (EPI) in the stomatogastric ganglion. A. For a choice of model (STG) and emergent property (network syncing), emergent property inference (EPI) learns a posterior distribution of the model parameters  $z = [g_{el}, g_{synA}]^\top$  conditioned on network syncing. B. An EPI distribution of STG model parameters producing network syncing. Samples are colored by log density. Distribution contours of emergent property value error are shown at levels of  $2 \times 10^{-6}$ ,  $2 \times 10^{-5}$ , and  $2 \times 10^{-4}$ . The eigenvectors of the Hessian at the mode of the inferred distribution are indicated as  $v_1$  and  $v_2$ . Simulated activity is shown for three samples (stars). (Inset) Sensitivity of the system with respect to network syncing along all dimensions of parameter space away from the mode. (see Section A.2.1). C. Deep probability distributions map a latent random variable  $w$  through a deep neural network with weights and biases  $\theta$  to parameters  $z$  distributed as  $q_\theta(z)$ . D. EPI learns a distribution  $q_\theta(z)$  of model parameters that produce an emergent property: the emergent property statistics  $T(x)$  are fixed in expectation over parameter distribution samples  $z \sim q_\theta(z)$  to particular values  $\mu$ .

121 Specifically, the two fast neurons ( $f_1$  and  $f_2$ ) mutually inhibit one another, and oscillate at a faster  
122 frequency than the mutually inhibiting slow neurons ( $s_1$  and  $s_2$ ), and the hub neuron (hub) couples  
123 with the fast or slow population or both.

124 Second, once the model is selected, one defines the emergent property, the measurable signal of  
125 scientific interest. To continue our running STG example, one such emergent property is the  
126 phenomenon of *network syncing* – in certain parameter regimes, the frequency of the hub neuron  
127 matches that of the fast and slow populations at an intermediate frequency. This emergent property  
128 is shown in Figure 1A at a frequency of 0.54Hz.

129 Third, qualitative parameter analysis ensues: since precise mathematical analysis is intractable in  
130 this model, a brute force sweep of parameters is done [22]. Subsequently, a qualitative description  
131 is formulated to describe the different parameter configurations that lead to the emergent property.  
132 In this last step lies the opportunity for a precise quantification of the emergent property as a  
133 statistical feature of the model. Once we have such a methodology, we can infer a probability  
134 distribution over parameter configurations that produce this emergent property.

135 Before presenting technical details (in the following section), let us understand emergent property  
136 inference schematically: EPI (Fig. 1A gray box) takes, as input, the model and the specified  
137 emergent property, and as its output, produces the parameter distribution shown in Figure 1B.  
138 This distribution – represented for clarity as samples from the distribution – is then a scientifically  
139 meaningful and mathematically tractable object. In the STG model, this distribution can be  
140 specifically queried to reveal the prototypical parameter configuration for network syncing (the  
141 mode; Figure 1B yellow star), and how network syncing decays based on changes away from the  
142 mode. Intuitively, the probability density of the samples is in agreement with the emergent property  
143 value error (Fig. 1B contours). Furthermore, the eigenvectors of the distribution Hessian at the  
144 mode can be queried to quantitatively formalize the robustness of network syncing (Fig. 1B  $v_1, v_2$ ).  
145 Indeed, samples equidistant from the mode along these EPI-identified dimensions of sensitivity ( $v_1$ )  
146 and degeneracy ( $v_2$ ) have diminished or preserved network syncing, respectively (Figure 1B inset  
147 and activity traces). Further validation of EPI is available in the supplementary materials, where  
148 we analyze a simpler model for which ground-truth statements can be made (Section A.1.1).

<sup>149</sup> **3.2 A deep generative modeling approach to emergent property inference**

<sup>150</sup> Emergent property inference (EPI) systematizes the three-step procedure of the previous section.  
<sup>151</sup> First, we consider the model as a coupled set of differential (and potentially stochastic) equations  
<sup>152</sup> [22]. In the running STG example, the dynamical state  $x = [x_{f1}, x_{f2}, x_{hub}, x_{s1}, x_{s2}]$  is the membrane  
<sup>153</sup> potential for each neuron, which evolves according to the biophysical conductance-based equation:

$$C_m \frac{dx}{dt} = -h(x; z) = -[h_{leak}(x; z) + h_{Ca}(x; z) + h_K(x; z) + h_{hyp}(x; z) + h_{elec}(x; z) + h_{syn}(x; z)] \quad (1)$$

<sup>154</sup> where  $C_m = 1\text{nF}$ , and  $h_{leak}$ ,  $h_{Ca}$ ,  $h_K$ ,  $h_{hyp}$ ,  $h_{elec}$ ,  $h_{syn}$  are the leak, calcium, potassium, hyperpolarization,  
<sup>155</sup> electrical, and synaptic currents, all of which have their own complicated dependence on  $x$   
<sup>156</sup> and  $z = [g_{el}, g_{synA}]$  (see Section A.2.1).

<sup>157</sup> Second, we define the emergent property, which as above is network syncing: oscillation of the  
<sup>158</sup> entire population at an intermediate frequency of our choosing (Figure 1A bottom). Quantifying  
<sup>159</sup> this phenomenon is straightforward: we define network syncing to be that each neuron’s spiking  
<sup>160</sup> frequency – denoted  $\omega_{f1}(x)$ ,  $\omega_{f2}(x)$ , etc. – is close to an intermediate frequency of 0.54Hz. Mathematically,  
<sup>161</sup> we achieve this via constraints on the mean and variance of  $\omega_i(x)$  for each neuron  
<sup>162</sup>  $i \in \{f1, f2, hub, s1, s2\}$ , and thus:

$$E[T(x)] \triangleq E \begin{bmatrix} \omega_{f1}(x) \\ \vdots \\ (\omega_{f1}(x) - 0.54)^2 \\ \vdots \end{bmatrix} = \begin{bmatrix} 0.54 \\ \vdots \\ 0.025^2 \\ \vdots \end{bmatrix} \triangleq \mu, \quad (2)$$

<sup>163</sup> which completes the quantification of the emergent property.

<sup>164</sup> Third, we perform emergent property inference: we find a distribution over parameter configura-  
<sup>165</sup> tions  $z$ , and insist that samples from this distribution produce the emergent property; in other  
<sup>166</sup> words, they obey the constraints introduced in Equation 2. This distribution will be chosen from  
<sup>167</sup> a family of probability distributions  $\mathcal{Q} = \{q_\theta(z) : \theta \in \Theta\}$ , defined by a deep generative distribution  
<sup>168</sup> of the normalizing flow class [16, 17, 18] – neural networks which transform a simple distribution  
<sup>169</sup> into a suitably complicated distribution (as is needed here). This deep distribution is represented  
<sup>170</sup> in Figure 1C (and see Methods for more detail). Then, mathematically, we must solve the following  
<sup>171</sup> optimization program:

$$\begin{aligned} & \underset{q_\theta \in \mathcal{Q}}{\operatorname{argmax}} H(q_\theta(z)) \\ & \text{s.t. } E_{z \sim q_\theta} [E_{x \sim p(x|z)} [T(x)]] = \mu, \end{aligned} \quad (3)$$

where  $T(x), \mu$  are defined as in Equation 2, and  $p(x|z)$  is the intractable distribution of data from the model ( $x$ ), given that model's parameters  $z$  (we access samples from this distribution by running the model forward). The purpose of each element in this program is detailed in Figure 1D. Finally, we recognize that many distributions in  $\mathcal{Q}$  will respect the emergent property constraints, so we require a normative principle to select amongst them. This principle is captured in Equation 3 by the primal objective  $H$ . Here we chose Shannon entropy as a means to find parameter distributions with minimal assumptions beyond some chosen structure [31, 32, 19, 33], but we emphasize that the EPI method is unaffected by this choice (but the results of course will depend on the primal objective chosen).

EPI optimizes the weights and biases  $\theta$  of the deep neural network (which induces the probability distribution) by iteratively solving Equation 3. The optimization is complete when the sampled models with parameters  $z \sim q_\theta$  produce activity consistent with the specified emergent property. Such convergence is evaluated with a hypothesis test that the mean of each emergent property statistic is not different than its emergent property value (see Section A.1.2). Equipped with this method, we now prove out the value of EPI by using it to investigate and produce novel insights about three prominent models in neuroscience.

### 3.3 Comprehensive input-responsivity in a nonlinear sensory system

Dynamical models of excitatory (E) and inhibitory (I) populations with supralinear input-output function have succeeded in explaining a host of experimentally documented phenomena. In a regime characterized by inhibitory stabilization of strong recurrent excitation, the model gives rise to paradoxical responses [4], selective amplification [34], surround suppression [35] and normalization [36]. Despite its strong predictive power, the E-I circuit model relies on the assumption that inhibition can be studied as an indivisible unit. Advances in experimental research reveal instead that inhibition is composed of distinct elements (parvalbumin (P), somatostatin(S), vip (V)) comprising 80% of GABAergic interneurons in V1 [37, 38, ?] and that these inhibitory cell types follow specific connectivity patterns (Fig. 2A) [56]. Recent theoretical advances [23, ?, ?], have only started to address the consequence of this multiplicity in the dynamics of V1, strongly relying on linear theory tools. Here, we use EPI to go beyond linear theory and systematically examine the distributions of parameters that are compatible with increases in neuron-type population rates, generating hypotheses of model operation.

Specifically, we consider a four-dimensional circuit model with dynamical state given by the firing

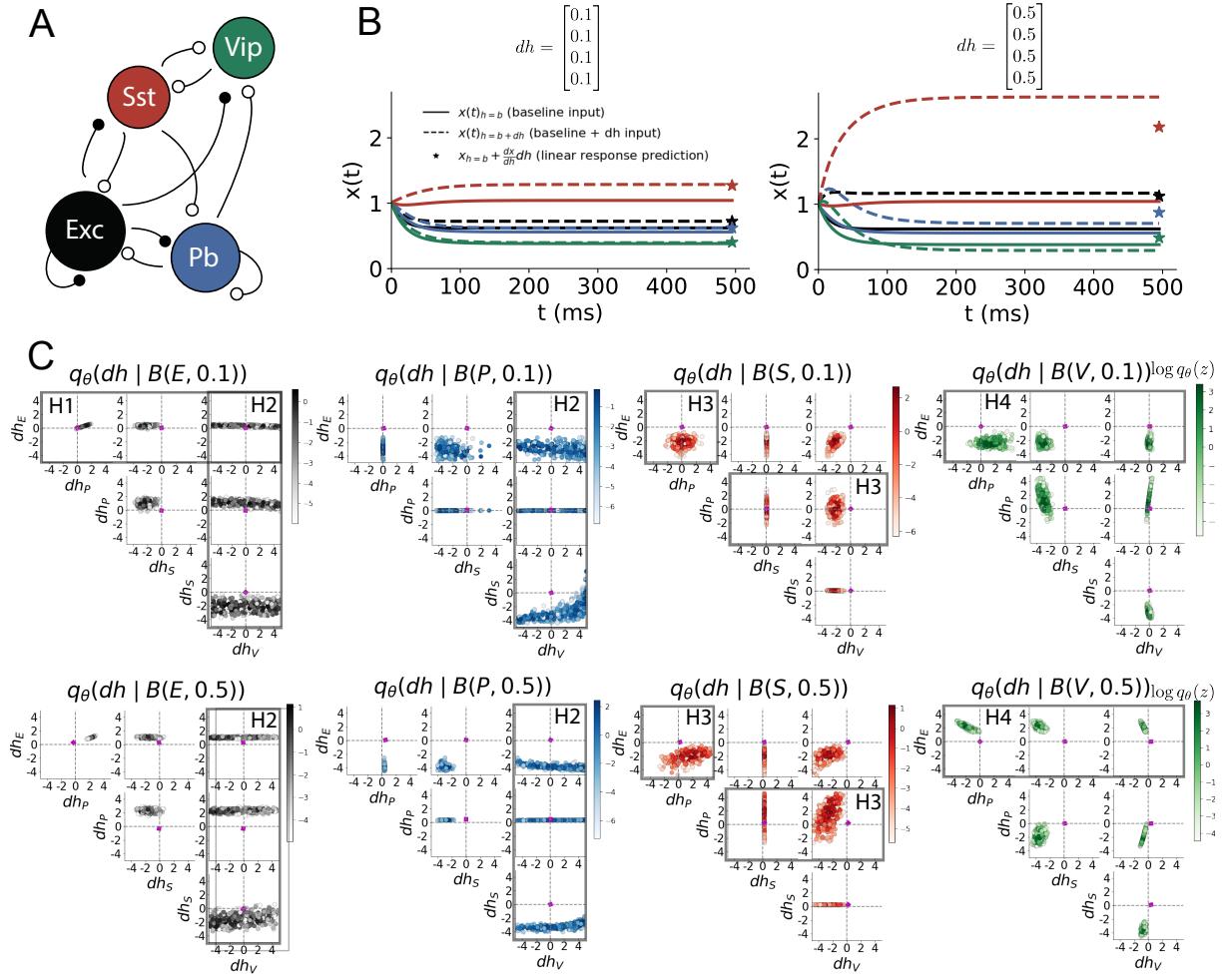


Figure 2: Hypothesis generation through EPI in a V1 model. A. Four-population model of primary visual cortex with excitatory (black), parvalbumin (blue), somatostatin (red), and vip (green) neurons. Some neuron-types largely do not form synaptic projections to others (excitatory and inhibitory projections filled and unfilled, respectively). B. Linear response predictions become inaccurate with greater input strength. V1 model simulations for input ( $h = b$ ) and ( $h = b + dh$ ) with  $b = [1, 1, 1, 1]^T$  and (left)  $dh = [0.1, 0.1, 0.1, 0.1]^T$  (right)  $dh = [0.5, 0.5, 0.5, 0.5]^T$ . Stars indicate the linear response prediction. C. EPI distributions on differential input  $dh$  conditioned on differential response  $\mathcal{B}(\alpha, y)$ . Supporting evidence for the four generated hypotheses are indicated by gray boxes with labels H1, H2, H3, and H4. The linear prediction from two standard deviations away from  $y$  (from negative to positive) is overlaid in magenta (very small, near origin).

rate  $x$  of each neuron-type population  $x = [x_E, x_P, x_S, x_V]^\top$ . Given a time constant of  $\tau = 20$  ms and a power  $n = 2$ , the dynamics are driven by the rectified ( $\|\cdot\|_+$ ) and exponentiated sum of recurrent ( $Wx$ ) and external  $h$  inputs:

$$\tau \frac{dx}{dt} = -x + [Wx + h]_+^n \quad (4)$$

The effective connectivity weights  $W$  were obtained from experimental recordings of publicly available datasets of mouse V1 [39, 40] (see Section A.2.2). The input  $h = b + dh$  is comprised of a baseline input  $b = [b_E, b_P, b_S, b_V]^\top$  and a differential input  $dh = [dh_E, dh_P, dh_S, dh_V]^\top$  to each neuron-type population. Throughout subsequent analyses, the baseline input is  $b = [1, 1, 1, 1]^\top$ .

With this model, we are interested in the differential responses of each neuron-type population to changes in input  $dh$ . Initially, we studied the linearized response of the system to input  $\frac{dx_{ss}}{dh}$  at the steady state response  $x_{ss}$ , i.e. a fixed point. All analyses of this model consider the steady state response, so we drop the notation  $ss$  from here on. While this linearization accurately predicts differential responses  $dx = [dx_E, dx_P, dx_S, dx_V]$  for small differential inputs to each population  $dh = [0.1, 0.1, 0.1, 0.1]$  (Fig 2B left), the linearization is a poor predictor in this nonlinear model more generally (Fig. 3B right). Currently available approaches to deriving the steady state response of the system are limited.

To get a more comprehensive picture of the input-responsivity of each neuron-type beyond linear theory, we used EPI to learn a distribution of the differential inputs to each population  $dh$  that produce an increase of  $y \in \{0.1, 0.5\}$  in the rate of each neuron-type population  $\alpha \in \{E, P, S, V\}$ . We want to know the differential inputs  $dh$  that result in a differential steady state  $dx_\alpha$  (the change in  $x_\alpha$  when receiving input  $h = b + dh$  with respect to the baseline  $h = b$ ) of value  $y$  with some small, arbitrarily chosen amount of variance  $0.01^2$ . These statements amount to the emergent property

$$\mathcal{B}(\alpha, y) \triangleq E \begin{bmatrix} dx_\alpha \\ (dx_\alpha - y)^2 \end{bmatrix} = \begin{bmatrix} y \\ 0.01^2 \end{bmatrix} \quad (5)$$

We maintain the notation  $\mathcal{B}(\cdot)$  throughout the rest of the study as short hand for emergent property, which represents a different signature of computation in each application. In each column of Figure 2C visualizes the inferred distribution of  $dh$  corresponding to a E (red), P (blue), S (red) and V (green) neuron-type increase, while each row corresponds to amounts of increase 0.1 and 0.5. These distributions conditioned on such emergent properties are now available through EPI. For each pair of parameters we show the two-dimensional marginal distribution of samples colored

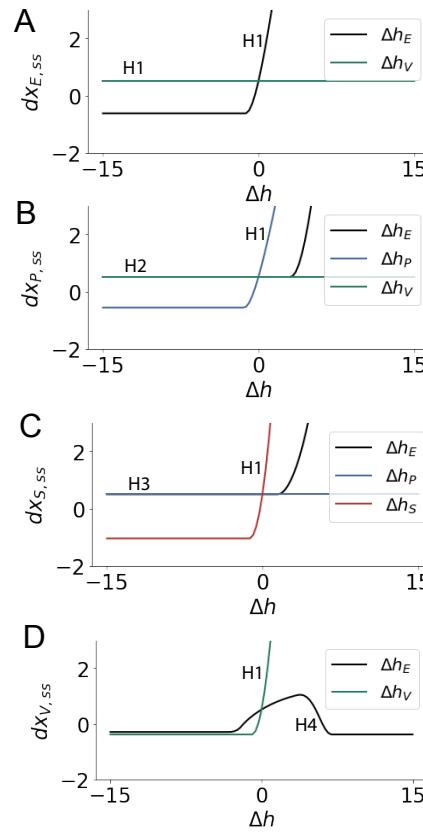


Figure 3: Confirming EPI generated hypotheses in V1. A. Differential responses by the E-population to changes in individual input  $\Delta h_\alpha u_\alpha$  away from the mode of the EPI distribution  $dh^*$ . B-D Same plots for the P-, S-, and V-populations. Labels H1, H2, H3, and H4 indicate which curves confirm which hypotheses.

230 by  $\log q_\theta(dh | \mathcal{B}(\alpha, y))$ . The inferred distributions immediately suggest four hypotheses:

231

- 232 H1: as is intuitive, each neuron-type's firing rate should be sensitive to that neuron-type's  
233 direct input (e.g. Fig. 2C H1 indicates low variance in  $dh_E$  when  $\alpha = E$ . Same observation  
234 in all inferred distributions);  
235 H2: the E- and P-populations should be largely unaffected by  $dh_V$  (Fig. 2C H2 indicates high variance in  $dh_V$  when  $\alpha \in \{E, P\}$ );  
236 H3: the S-population should be largely unaffected by  $dh_P$  (Fig. 2C H3 indicate high variance in  $dh_P$  when  $\alpha = S$ );  
237 H4: there should be a nonmonotonic response of  $dx_{V,ss}$  with  $dh_E$  (Fig. 2C H4 indicates that  
238 negative  $dh_E$  should result in small  $dx_{V,ss}$ , but positive  $dh_E$  should elicit a larger  $dx_{V,ss}$ );

239  
240 We evaluate these hypotheses by taking steps in individual neuron-type input  $\Delta h_\alpha$  away from the  
241 modes of the inferred distributions at  $y = 0.1$ .

$$dh^* = z^* = \underset{z}{\operatorname{argmax}} \log q_\theta(z | \mathcal{B}(\alpha, 0.1)) \quad (6)$$

242 Now,  $\Delta x_\alpha$  is the change in steady state response to the system with input  $h = b + dh^* + \Delta h_\alpha u_\alpha$

244 compared to  $h = b + dh^*$ , where  $u_\alpha$  is a unit vector in the dimension of  $\alpha$ . The EPI-generated  
 245 hypotheses are confirmed.

246 H1: the neuron-type responses are sensitive to their direct inputs (Fig. 3A black, 3B blue,  
 247 3C red, 3D green);

248 H2: the E- and P-populations are not affected by  $dh_V$  (Fig. 3A green, 3B green);

249 H3: the S-population is not affected by  $dh_P$  (Fig. 3C blue);

250 H4: the V-population exhibits a nonmonotonic response to  $dh_E$  (Fig. 3D black), and is in  
 251 fact the on population to do so (Fig. 3A-C black).

252 These hypotheses were in stark contrast to what was available to us via traditional analytical linear  
 253 prediction (Fig. 2C, magenta). To this point, we have shown the utility of EPI on relatively low-  
 254 level emergent properties like network syncing and differential neuron-type population responses.  
 255 In the remainder of the study, we focus on using EPI to understand models of more abstract  
 256 cognitive function.

### 257 3.4 Identifying neural mechanisms of behavioral learning.

258 Identifying measurable biological changes that result in improved behavior is important for neuro-  
 259 science, since they may indicate how the learning brain adapts. In a rapid task switching experiment  
 260 [41], rats were explicitly cued on each trial to either orient towards a visual stimulus in the Pro  
 261 (P) task or orient away from a visual stimulus in the Anti (A) task (Fig. 3a). Neural recordings  
 262 in the midbrain supeior colliculus (SC) exhibited two population of neurons that simultaneously  
 263 represented both task context (Pro or Anti) and motor response (contralateral or ipsilateral to the  
 264 recorded side): the Pro/Contra and Anti/Ipsi neurons [24]. Duan et al. proposed a model of SC  
 265 that, like the V1 model analyzed in the previous section, is a four-population dynamical system.  
 266 Here, the neuron-type populations are functionally-defined as the Pro- and Anti-populations in each  
 267 hemisphere (left (L) and right (R)). The Pro- or Anti-populations receive an input determined by  
 268 the cue, and then the left and right populations receive an input based on the side of the light  
 269 stimulus. Activities were bounded between 0 and 1, so that a high output of the Pro population  
 270 in a given hemisphere corresponds to the contralateral response. An additional stipulation is that  
 271 when one Pro population responds with a high-output, the opposite Pro population must respond  
 272 with a low output. Finally, this circuit operates in the presence of Gaussian noise resulting in trial-  
 273 to-trial variability (see Section A.2.3). The connectivity matrix is parameterized by the geometry  
 274 of the population arrangement (Fig. 3B).

275 Here, we used EPI to learn distributions of the SC weight matrix parameters  $z = W$  conditioned  
 276 on of various levels of rapid task switching accuracy  $\mathcal{B}(p)$  for  $p \in \{50\%, 60\%, 70\%, 80\%, 90\%\}$  (see  
 277 Section A.2.3). Following the approach in Duan et al., we decomposed the connectivity matrix  
 278  $W = QAQ^{-1}$  in such a way (the Schur decomposition) that the basis vectors  $q_i$  are the same for all  
 279  $W$  (Fig. 3C). These basis vectors have intuitive roles in processing for this task, and are accordingly  
 280 named the *all* mode - all neurons co-fluctuate, *side* mode - one side dominates the other, *task* mode  
 281 - the Pro or Anti populations dominate the other, and *diag* mode - Pro- and Anti-populations of  
 282 opposite hemispheres dominate the opposite pair. The corresponding eigenvalues (e.g.  $a_{\text{task}}$ , which  
 283 change according to  $W$ ) indicate the degree to which activity along that mode is increased or  
 284 decreased by  $W$ .

285 EPI demonstrates that, for greater task accuracies, the task mode eigenvalue increases, indicating  
 286 the importance of  $W$  to the task representation (Fig. 4D, purple). Stepping from random chance  
 287 (50%) networks to marginally task-performing (60%) networks, there is a marked decrease of the  
 288 side mode eigenvalues (Fig. 3D, orange). Such side mode suppression remains in the models  
 289 achieving greater accuracy, revealing its importance towards task performance. There were no  
 290 interesting trends with learning in the all or diag mode (hence not shown in Fig. 3). Importantly,  
 291 we can conclude from our methodology that side mode suppression in  $W$  allows rapid task switching,  
 292 and that greater task-mode representations in  $W$  increase accuracy. These hypotheses are confirmed  
 293 by forward simulation of the SC model (Fig. 3E). Thus, EPI produces novel, experimentally testable  
 294 predictions: effective connectivity between these populations changes throughout learning, in a way  
 295 that increases its task mode and decreases its side mode eigenvalues.

### 296 3.5 Linking RNN connectivity to computational error

297 So far, each model we have studied was designed from fundamental biophysical principles, genetically-  
 298 or functionally-defined neuron types. At a more abstract level of modeling, recurrent neural net-  
 299 works (RNNs) are high-dimensional dynamical models of computation that are becoming increas-  
 300 ingly popular in neuroscience research [42]. In theoretical neuroscience, RNN dynamics usually  
 301 follow the equation

$$\frac{dx}{dt} = -x(t) + W\phi(x(t)) + I(t), \quad (7)$$

302 where  $x(t)$  is the network activity,  $W$  is the network connectivity,  $\phi(\cdot) = \tanh(\cdot)$ , and  $I(t)$  is the  
 303 input to the system. Such RNNs are trained to do a task from a systems neuroscience experiment,  
 304 and then the unit activations of the trained RNN are compared to recorded neural activity. Fully-

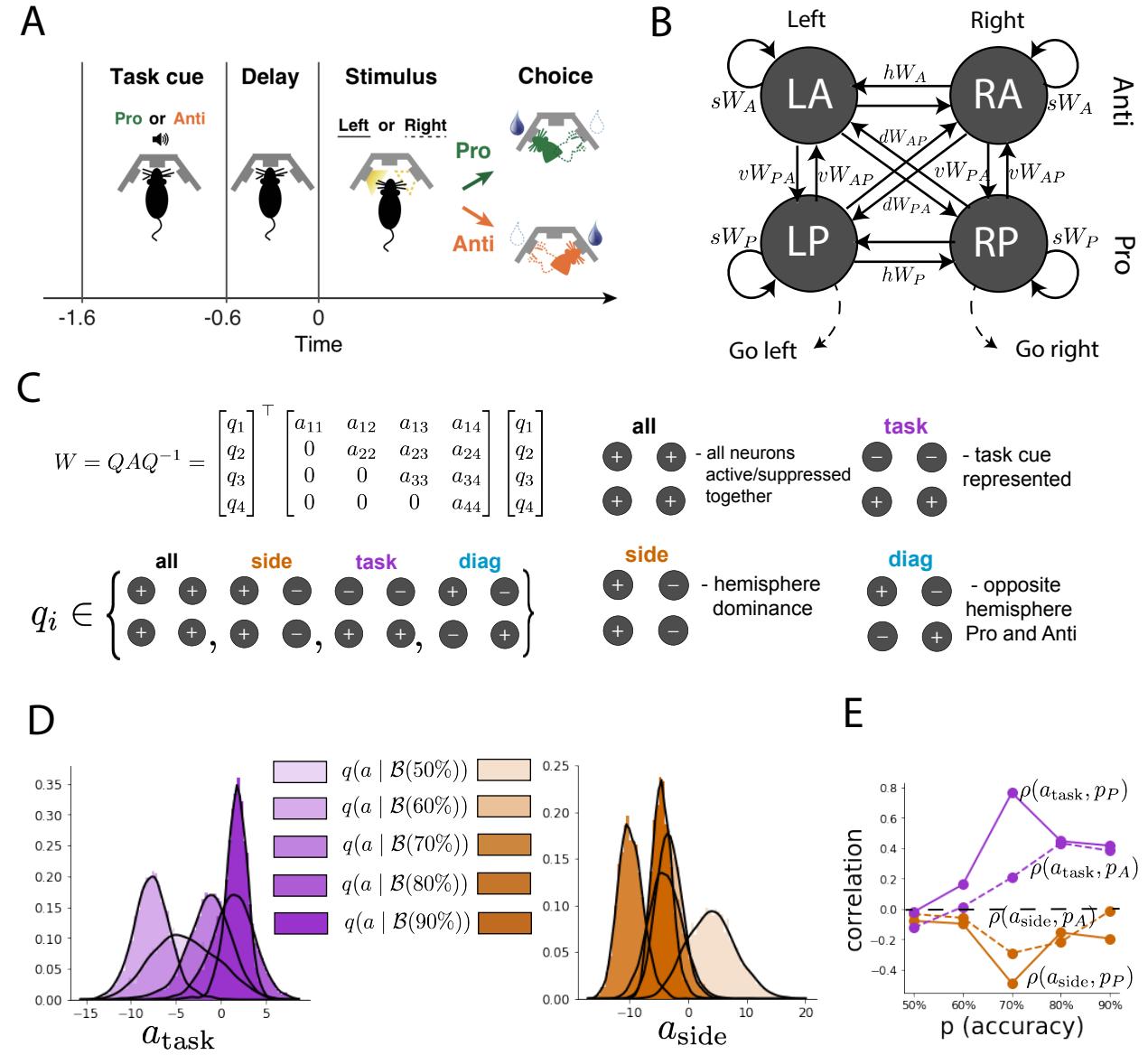


Figure 4: EPI reveals changes in SC [24] connectivity that control task accuracy. A. Rapid task switching behavioral paradigm (see text). B. Model of superior colliculus (SC). Neurons: LP - left pro, RP - right pro, LA - left anti, RA - right anti. Parameters:  $sW$  - self,  $hW$  - horizontal,  $vW$  - vertical,  $dW$  - diagonal weights. C. The Schur decomposition of the weight matrix  $W = QAQ^{-1}$  is a unique decomposition with orthogonal  $Q$  and upper triangular  $A$ . Schur modes:  $q_{all}$ ,  $q_{task}$ ,  $q_{side}$ , and  $q_{diag}$ . D. The marginal EPI distributions of the Schur eigenvalues at each level of task accuracy. E. The correlation of Schur eigenvalue with task performance in each learned EPI distribution.

305 connected RNNs with tens of thousands of parameters are challenging to characterize [43], especially  
 306 making statistical inferences about their parameterization. Alternatively, we consider a rank-1,  $N$ -  
 307 neuron RNN with connectivity

$$W = g\chi + \frac{1}{N}mn^\top, \quad (8)$$

308 where  $\chi_{ij} \sim \mathcal{N}(0, \frac{1}{N})$ ,  $g$  is the random strength, and the entries of  $m$  and  $n$  are drawn from Gaussian  
 309 distributions  $m_i \sim \mathcal{N}(M_m, 1)$  and  $n_i \sim \mathcal{N}(M_n, 1)$ . We use EPI to infer the parameterizations of  
 310 rank-1 RNNs solving an example task, enabling discovery of properties of connectivity that result  
 311 in different types of computational errors.

312 The task we consider is Gaussian posterior conditioning: calculate the parameters of a posterior  
 313 distribution induced by a prior  $p(\mu_y) = \mathcal{N}(\mu_0 = 4, \sigma_0^2 = 1)$  and a likelihood  $p(y|\mu_y) = \mathcal{N}(\mu_y, \sigma_y^2 =$   
 314 1), given a single observation  $y$ . Conjugacy offers the result analytically;  $p(\mu_y|y) = \mathcal{N}(\mu_{post}, \sigma_{post}^2)$ ,  
 315 where:

$$\mu_{post} = \frac{\frac{\mu_0}{\sigma_0^2} + \frac{y}{\sigma_y^2}}{\frac{1}{\sigma_0^2} + \frac{1}{\sigma_y^2}} \quad \sigma_{post}^2 = \frac{1}{\frac{1}{\sigma_0^2} + \frac{1}{\sigma_y^2}}. \quad (9)$$

316 The RNN is trained to solve this task by producing readout activity that is on average the posterior  
 317 mean  $\mu_{post}$ , and activity whose variability is the posterior variance  $\sigma_{post}^2$  (a setup inspired by  
 318 [44]). To solve this Gaussian posterior conditioning task, the RNN response to a constant input  
 319  $I(t) = yw + (n - M_n)$  must equal the posterior mean along readout vector  $w$ , where

$$\kappa_w = \frac{1}{N} \sum_{j=1}^N w_j \phi(x_j) \quad (10)$$

320 Additionally, the amount of chaotic variance  $\Delta_T$  must equal the posterior variance.  $\kappa_w$  and  $\Delta_T$  can  
 321 be expressed in terms of each other through a solvable system of nonlinear equations (see Section  
 322 A.2.4) [25]. This theory allows us to mathematically formalize the execution of this task into an  
 323 emergent property, where the emergent property statistics of the RNN activity are  $k_w$  and  $\Delta_T$  and  
 324 the emergent property values are the ground truth posterior mean  $\mu_{post}$  and variance  $\sigma_{post}^2$ :

$$E \begin{bmatrix} \kappa_w \\ \Delta_T \\ (\kappa_w - \mu_{post})^2 \\ (\Delta_T - \sigma_{post}^2)^2 \end{bmatrix} = \begin{bmatrix} \mu_{post} \\ \sigma_{post}^2 \\ 0.1 \\ 0.1 \end{bmatrix} \quad (11)$$

325 We specify a substantial amount of variability in the variance constraints so that the inferred  
 326 distribution results in RNNs with a variety biases in their solutions to the gaussian posterior  
 327 conditioning problem.

328 We used EPI to learn distributions of RNN connectivity properties  $z = [g \ M_m \ M_n]$  executing  
 329 Gaussian posterior conditioning given an input of  $y = 2$ . (see Section A.2.4) (Fig. 5B). The true  
 330 Gaussian conditioning posterior for an input of  $y = 2$  is  $\mu_{\text{post}} = 3$  and  $\sigma_{\text{post}} = 0.5$ . We examined  
 331 the nature of the over- and under-estimation of the posterior means (Fig. 5B, left) and variances  
 332 (Fig. 5B, right) in the inferred distributions. There is rough symmetry in the  $M_m$ - $M_n$  plane,  
 333 suggesting a degeneracy in the product of  $M_m$  and  $M_n$  (Fig. 5B). The product of  $M_m$  and  $M_n$   
 334 almost completely determines the posterior mean (Fig. 5B, left), and the random strength  $g$  is the  
 335 most influential variable on the temporal variance (Fig. 5B, right). Neither of these observations  
 336 were obvious from what mathematical analysis is available in networks of this type (see Section  
 337 A.2.4). They lead to the following hypotheses:

- 338 H1: The posterior mean of the RNN increases with the product of  $M_m$  and  $M_n$ ;  
 339 H2: The posterior variance increases with  $g$ ;

340

341 Testing these now in finite-size networks. Will write end of this later.  
 342 This novel procedure of doing inference in interpretable parameterizations of RNNs conditioned on  
 343 the emergent property of task execution is straightforwardly generalizable to other tasks like noisy  
 344 integration and context-dependent decision making (Fig. S1).

345 **4 Discussion**

346 **4.1 EPI is a general tool for theoretical neuroscience**

347 Models of biological systems are often comprised of complex nonlinear differential equations, mak-  
 348 ing traditional theoretical analysis and statistical inference intractable. In contrast, EPI is capable  
 349 of learning distributions of parameters in such models producing measurable signatures of compu-  
 350 tation. We have demonstrated its utility on biological models (STG), intermediate-level models of  
 351 interacting genetically- and functionally-defined neuron-types (V1, SC), and the most abstract of  
 352 models (RNNs). We are able to condition both deterministic and stochastic models on low-level  
 353 emergent properties like firing rates of membrane potentials, as well as high-level cognitive func-  
 354 tion like Gaussian posterior conditioning. Technically, EPI is tractable when the emergent property  
 355 statistics are continuously differentiable with respect to the model parameters, which is very often  
 356 the case; this emphasizes the general utility of EPI.

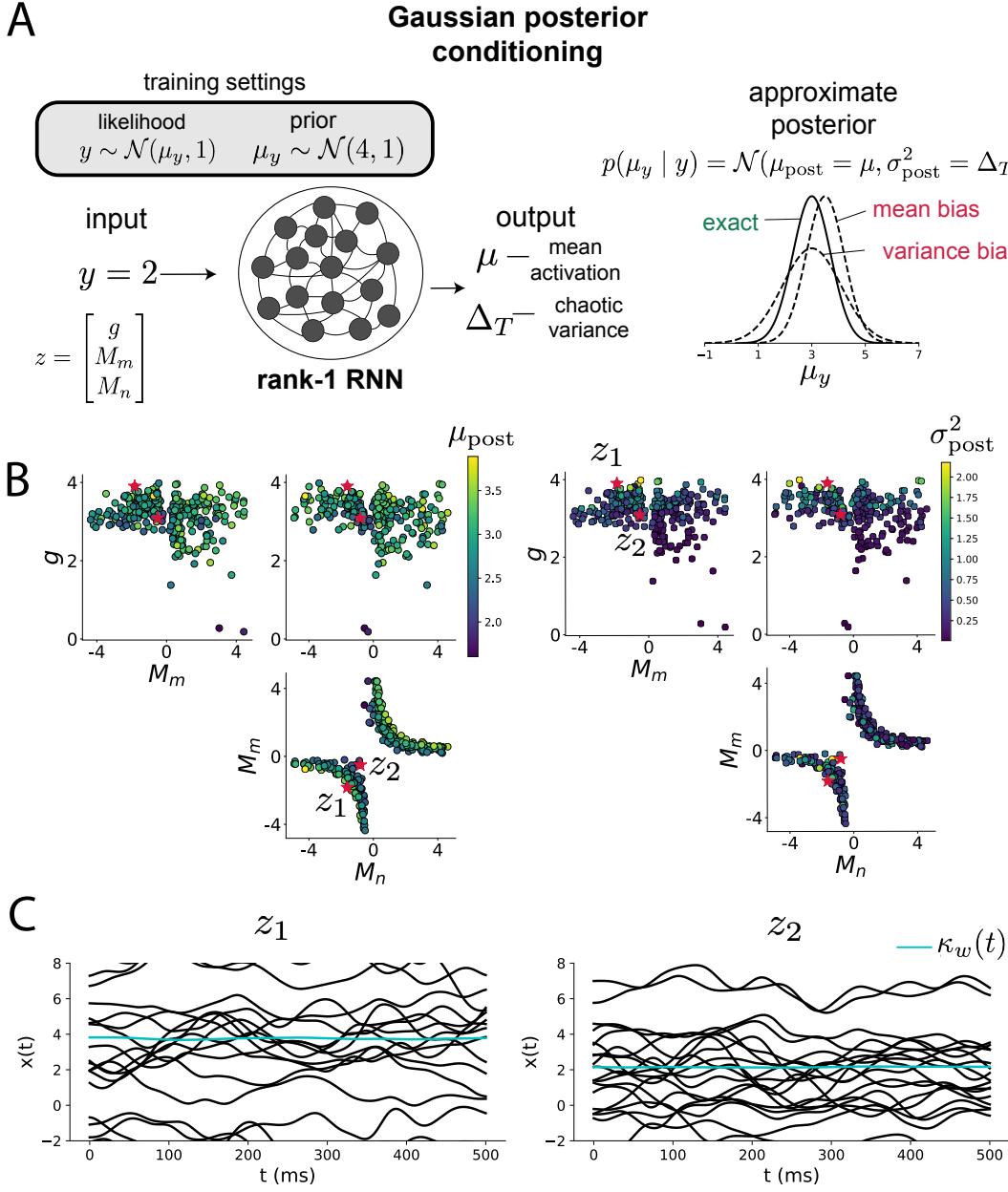


Figure 5: Sources of solution bias in an RNN computation. A. (left) A rank-1 RNN executing a Gaussian posterior conditioning computation on  $\mu_y$ . (right) Bias in this computation can come from over- or under-estimating the posterior mean or variance. B. EPI distribution of rank-1 RNNs executing Gaussian posterior conditioning. Samples are colored by (left) posterior mean  $\mu_{\text{post}} = \kappa_w$  and (right) posterior variance  $\sigma_{\text{post}}^2 = \Delta_T$ . C. Finite-size networks sampled from the distribution perform the calculation and have the computational biases expected from their parameter values. Activity along readout  $\kappa_w$  (cyan).

357 In this study, we have focused on applying EPI to low dimensional parameter spaces of models  
358 with low dimensional dynamical state. These choices were made to present the reader with a series  
359 of interpretable conclusions, which is more challenging in high dimensional spaces. In fact, EPI  
360 should scale reasonably to high dimensional parameter spaces, as the underlying technology has  
361 produced state-of-the-art performance on high-dimensional tasks such as texture generation [19].  
362 Of course, increasing the dimensionality of the dynamical state of the model makes optimization  
363 more expensive, and there is a practical limit there as with any machine learning approach. For  
364 systems with high dimensional state, we recommend using theoretical approaches (e.g. [25]) to  
365 reason about reduced parameterizations of such high-dimensional systems.  
366 There are additional technical considerations when assessing the suitability of EPI for a particu-  
367 lar modeling question. First and foremost, as in any optimization problem, the defined emergent  
368 property should always be appropriately conditioned (constraints should not have wildly different  
369 units). Furthermore, if the program is underconstrained (not enough constraints), the distribution  
370 grows (in entropy) unstably unless mapped to a finite support. If overconstrained, there is no pa-  
371 rameter set producing the emergent property, and EPI optimization will fail (appropriately). Next,  
372 one should consider the computational cost of the gradient calculations. In the best circumstance,  
373 there is a simple, closed form expression (e.g. Section A.1.1) for the emergent property statistic  
374 given the model parameters. On the other end of the spectrum, many forward simulation iterations  
375 may be required before a high quality measurement of the emergent property statistic is available  
376 (e.g. Section A.2.1). In such cases, optimization will be expensive.

## 377 4.2 Novel hypotheses from EPI

378 Machine learning has played an effective, multifaceted role in neuroscientific progress. Primarily,  
379 it has revealed structure in large-scale neural datasets [45, 46, 47, 48, 49, 50] (see review, [14]).  
380 Secondarily, trained algorithms of varying degrees of biological relevance are beginning to be viewed  
381 as fully-observable computational systems comparable to the brain [43, 51].  
382 For example, consider the fact that we do not fully understand the four-dimensional models of V1  
383 [23]. Because analytical approaches to studying nonlinear dynamical systems become increasingly  
384 complicated when stepping from two-dimensional to three- or four-dimensional systems in the  
385 absence of restrictive simplifying assumptions [52], it is unsurprising that this model has been a  
386 challenge. In Section 3.3, we showed that EPI was far more informative about neuron-type input  
387 responsibility than the predictions afforded through analysis. By flexibly conditioning this V1 model

388 on different emergent properties, we performed an exploratory analysis of a *model* rather than a  
389 dataset, which generated and proved out a set of testable predictions.

390 Of course, exploratory analyses can also be directed. For example, when interested in model  
391 changes during learning, one can use EPI to condition as we did in Section 3.4. This analysis  
392 identified experimentally testable predictions (proved out *in-silico*) of changes in connectivity in  
393 SC throughout learning. Precisely, we predict that an initial reduction in side mode eigenvalue,  
394 and a steady increase in task mode eigenvalue will take place, during learning, in the effective  
395 connectivity matrices of learning rats.

396 In our final analysis, we present a novel procedure for doing statistical inference on interpretable  
397 parameterizations of RNNs executing simple tasks . This methodology relies on recently extended  
398 theory of responses in random neural networks with minimal structure [25]. With this methodology,  
399 we can finally open the probabilistic model selection toolkit reasoning about the connectivity of  
400 RNNs solving tasks.

## 401 References

- 402 [1] Larry F Abbott. Theoretical neuroscience rising. *Neuron*, 60(3):489–495, 2008.
- 403 [2] John J Hopfield. Neural networks and physical systems with emergent collective computational  
404 abilities. *Proceedings of the national academy of sciences*, 79(8):2554–2558, 1982.
- 405 [3] Haim Sompolinsky, Andrea Crisanti, and Hans-Jurgen Sommers. Chaos in random neural  
406 networks. *Physical review letters*, 61(3):259, 1988.
- 407 [4] Misha V Tsodyks, William E Skaggs, Terrence J Sejnowski, and Bruce L McNaughton. Para-  
408 doxical effects of external modulation of inhibitory interneurons. *Journal of neuroscience*,  
409 17(11):4382–4388, 1997.
- 410 [5] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *International Confer-  
411 ence on Learning Representations*, 2014.
- 412 [6] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation  
413 and variational inference in deep latent gaussian models. *International Conference on Machine  
414 Learning*, 2014.

- [7] Yuanjun Gao, Evan W Archer, Liam Paninski, and John P Cunningham. Linear dynamical neural population models through nonlinear embeddings. In *Advances in neural information processing systems*, pages 163–171, 2016.
- [8] Yuan Zhao and Il Memming Park. Recursive variational bayesian dual estimation for nonlinear dynamics and non-gaussian observations. *stat*, 1050:27, 2017.
- [9] Gabriel Barello, Adam Charles, and Jonathan Pillow. Sparse-coding variational auto-encoders. *bioRxiv*, page 399246, 2018.
- [10] Chethan Pandarinath, Daniel J O’Shea, Jasmine Collins, Rafal Jozefowicz, Sergey D Stavisky, Jonathan C Kao, Eric M Trautmann, Matthew T Kaufman, Stephen I Ryu, Leigh R Hochberg, et al. Inferring single-trial neural population dynamics using sequential auto-encoders. *Nature methods*, page 1, 2018.
- [11] Alexander B Wiltschko, Matthew J Johnson, Giuliano Iurilli, Ralph E Peterson, Jesse M Katon, Stan L Pashkovski, Victoria E Abraira, Ryan P Adams, and Sandeep Robert Datta. Mapping sub-second structure in mouse behavior. *Neuron*, 88(6):1121–1135, 2015.
- [12] Matthew J Johnson, David K Duvenaud, Alex Wiltschko, Ryan P Adams, and Sandeep R Datta. Composing graphical models with neural networks for structured representations and fast inference. In *Advances in neural information processing systems*, pages 2946–2954, 2016.
- [13] Eleanor Batty, Matthew Whiteway, Shreya Saxena, Dan Biderman, Taiga Abe, Simon Musall, Winthrop Gillis, Jeffrey Markowitz, Anne Churchland, John Cunningham, et al. Behavenet: nonlinear embedding and bayesian neural decoding of behavioral videos. *Advances in Neural Information Processing Systems*, 2019.
- [14] Liam Paninski and John P Cunningham. Neural data science: accelerating the experiment-analysis-theory cycle in large-scale neuroscience. *Current opinion in neurobiology*, 50:232–241, 2018.
- [15] Mark K Transtrum, Benjamin B Machta, Kevin S Brown, Bryan C Daniels, Christopher R Myers, and James P Sethna. Perspective: Sloppiness and emergent theories in physics, biology, and beyond. *The Journal of chemical physics*, 143(1):07B201\_1, 2015.
- [16] Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. *International Conference on Machine Learning*, 2015.

- 444 [17] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp.  
445 *arXiv preprint arXiv:1605.08803*, 2016.
- 446 [18] George Papamakarios, Theo Pavlakou, and Iain Murray. Masked autoregressive flow for density  
447 estimation. In *Advances in Neural Information Processing Systems*, pages 2338–2347, 2017.
- 448 [19] Gabriel Loaiza-Ganem, Yuanjun Gao, and John P Cunningham. Maximum entropy flow  
449 networks. *International Conference on Learning Representations*, 2017.
- 450 [20] Dustin Tran, Rajesh Ranganath, and David Blei. Hierarchical implicit models and likelihood-  
451 free variational inference. In *Advances in Neural Information Processing Systems*, pages 5523–  
452 5533, 2017.
- 453 [21] Mark S Goldman, Jorge Golowasch, Eve Marder, and LF Abbott. Global structure, robustness,  
454 and modulation of neuronal models. *Journal of Neuroscience*, 21(14):5229–5238, 2001.
- 455 [22] Gabrielle J Gutierrez, Timothy O’Leary, and Eve Marder. Multiple mechanisms switch an  
456 electrically coupled, synaptically inhibited neuron between competing rhythmic oscillators.  
457 *Neuron*, 77(5):845–858, 2013.
- 458 [23] Ashok Litwin-Kumar, Robert Rosenbaum, and Brent Doiron. Inhibitory stabilization and vi-  
459 sual coding in cortical circuits with multiple interneuron subtypes. *Journal of neurophysiology*,  
460 115(3):1399–1409, 2016.
- 461 [24] Chunyu A Duan, Marino Pagan, Alex T Piet, Charles D Kopec, Athena Akrami, Alexander J  
462 Riordan, Jeffrey C Erlich, and Carlos D Brody. Collicular circuits for flexible sensorimotor  
463 routing. *bioRxiv*, page 245613, 2018.
- 464 [25] Francesca Mastrogiovanni and Srdjan Ostojic. Linking connectivity, dynamics, and computa-  
465 tions in low-rank recurrent neural networks. *Neuron*, 99(3):609–623, 2018.
- 466 [26] Sean R Bittner, Agostina Palmigiano, Kenneth D Miller, and John P Cunningham. Degener-  
467 ate solution networks for theoretical neuroscience. *Computational and Systems Neuroscience  
468 Meeting (COSYNE), Lisbon, Portugal*, 2019.
- 469 [27] Sean R Bittner, Alex T Piet, Chunyu A Duan, Agostina Palmigiano, Kenneth D Miller,  
470 Carlos D Brody, and John P Cunningham. Examining models in theoretical neuroscience with  
471 degenerate solution networks. *Bernstein Conference*, 2019.

- 472 [28] Jan-Matthis Lueckmann, Pedro Goncalves, Chaitanya Chintaluri, William F Podlaski, Giacomo Bassetto, Tim P Vogels, and Jakob H Macke. Amortised inference for mechanistic models  
473 of neural dynamics. In *Computational and Systems Neuroscience Meeting (COSYNE), Lisbon, Portugal*, 2019.
- 476 [29] Eve Marder and Vatsala Thirumalai. Cellular, synaptic and network effects of neuromodulation.  
477 *Neural Networks*, 15(4-6):479–493, 2002.
- 478 [30] Astrid A Prinz, Dirk Bucher, and Eve Marder. Similar network activity from disparate circuit  
479 parameters. *Nature neuroscience*, 7(12):1345, 2004.
- 480 [31] Edwin T Jaynes. Information theory and statistical mechanics. *Physical review*, 106(4):620,  
481 1957.
- 482 [32] Gamaleldin F Elsayed and John P Cunningham. Structure in neural population recordings:  
483 an expected byproduct of simpler phenomena? *Nature neuroscience*, 20(9):1310, 2017.
- 484 [33] Cristina Savin and Gašper Tkačik. Maximum entropy models as a tool for building precise  
485 neural controls. *Current opinion in neurobiology*, 46:120–126, 2017.
- 486 [34] Brendan K Murphy and Kenneth D Miller. Balanced amplification: a new mechanism of  
487 selective amplification of neural activity patterns. *Neuron*, 61(4):635–648, 2009.
- 488 [35] Hirofumi Ozeki, Ian M Finn, Evan S Schaffer, Kenneth D Miller, and David Ferster. Inhibitory  
489 stabilization of the cortical network underlies visual surround suppression. *Neuron*, 62(4):578–  
490 592, 2009.
- 491 [36] Daniel B Rubin, Stephen D Van Hooser, and Kenneth D Miller. The stabilized supralinear  
492 network: a unifying circuit motif underlying multi-input integration in sensory cortex. *Neuron*,  
493 85(2):402–417, 2015.
- 494 [37] Henry Markram, Maria Toledo-Rodriguez, Yun Wang, Anirudh Gupta, Gilad Silberberg, and  
495 Caizhi Wu. Interneurons of the neocortical inhibitory system. *Nature reviews neuroscience*,  
496 5(10):793, 2004.
- 497 [38] Bernardo Rudy, Gordon Fishell, SooHyun Lee, and Jens Hjerling-Leffler. Three groups of  
498 interneurons account for nearly 100% of neocortical gabaergic neurons. *Developmental neuro-  
499 biology*, 71(1):45–61, 2011.

- 500 [39] (2018) Allen Institute for Brain Science. Layer 4 model of v1. available from:  
501 <https://portal.brain-map.org/explore/models/l4-mv1>.
- 502 [40] Yanzan N Billeh, Binghuang Cai, Sergey L Gratiy, Kael Dai, Ramakrishnan Iyer, Nathan W  
503 Gouwens, Reza Abbasi-Asl, Xiaoxuan Jia, Joshua H Siegle, Shawn R Olsen, et al. Systematic  
504 integration of structural and functional data into multi-scale models of mouse primary visual  
505 cortex. *bioRxiv*, page 662189, 2019.
- 506 [41] Chunyu A Duan, Jeffrey C Erlich, and Carlos D Brody. Requirement of prefrontal and midbrain  
507 regions for rapid executive control of behavior in the rat. *Neuron*, 86(6):1491–1503, 2015.
- 508 [42] Omri Barak. Recurrent neural networks as versatile tools of neuroscience research. *Current  
509 opinion in neurobiology*, 46:1–6, 2017.
- 510 [43] David Sussillo and Omri Barak. Opening the black box: low-dimensional dynamics in high-  
511 dimensional recurrent neural networks. *Neural computation*, 25(3):626–649, 2013.
- 512 [44] Rodrigo Echeveste, Laurence Aitchison, Guillaume Hennequin, and Máté Lengyel. Cortical-like  
513 dynamics in recurrent circuits optimized for sampling-based probabilistic inference. *bioRxiv*,  
514 page 696088, 2019.
- 515 [45] Robert E Kass and Valérie Ventura. A spike-train probability model. *Neural computation*,  
516 13(8):1713–1720, 2001.
- 517 [46] Emery N Brown, Loren M Frank, Dengda Tang, Michael C Quirk, and Matthew A Wilson.  
518 A statistical paradigm for neural spike train decoding applied to position prediction from  
519 ensemble firing patterns of rat hippocampal place cells. *Journal of Neuroscience*, 18(18):7411–  
520 7425, 1998.
- 521 [47] Liam Paninski. Maximum likelihood estimation of cascade point-process neural encoding  
522 models. *Network: Computation in Neural Systems*, 15(4):243–262, 2004.
- 523 [48] M Yu Byron, John P Cunningham, Gopal Santhanam, Stephen I Ryu, Krishna V Shenoy, and  
524 Maneesh Sahani. Gaussian-process factor analysis for low-dimensional single-trial analysis  
525 of neural population activity. In *Advances in neural information processing systems*, pages  
526 1881–1888, 2009.

- 527 [49] Kenneth W Latimer, Jacob L Yates, Miriam LR Meister, Alexander C Huk, and Jonathan W  
 528 Pillow. Single-trial spike trains in parietal cortex reveal discrete steps during decision-making.  
 529 *Science*, 349(6244):184–187, 2015.
- 530 [50] Lea Duncker, Gergo Bohner, Julien Boussard, and Maneesh Sahani. Learning interpretable  
 531 continuous-time models of latent stochastic dynamical systems. *Proceedings of the 36th Inter-*  
 532 *national Conference on Machine Learning*, 2019.
- 533 [51] Blake A Richards and et al. A deep learning framework for neuroscience. *Nature Neuroscience*,  
 534 2019.
- 535 [52] Steven H Strogatz. Nonlinear dynamics and chaos: with applications to physics. *Biology,*  
 536 *Chemistry, and Engineering (Studies in Nonlinearity)*, Perseus, Cambridge, UK, 1994.
- 537 [53] Rajesh Ranganath, Sean Gerrish, and David Blei. Black box variational inference. In *Artificial*  
 538 *Intelligence and Statistics*, pages 814–822, 2014.
- 539 [54] Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and  
 540 variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.
- 541 [55] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp.  
 542 *Proceedings of the 5th International Conference on Learning Representations*, 2017.
- 543 [56] Carsten K Pfeffer, Mingshan Xue, Miao He, Z Josh Huang, and Massimo Scanziani. Inhi-  
 544 bition of inhibition in visual cortex: the logic of connections between molecularly distinct  
 545 interneurons. *Nature Neuroscience*, 16(8):1068, 2013.

546 **A Methods**

547 **A.1 Emergent property inference (EPI)**

548 Emergent property inference (EPI) learns distributions of theoretical model parameters that pro-  
 549 duce emergent properties of interest. EPI combines ideas from likelihood-free variational inference  
 550 [20] and maximum entropy flow networks [19]. A maximum entropy flow network is used as a deep  
 551 probability distribution for the parameters, while these samples often parameterize a differentiable  
 552 model simulator, which may lack a tractable likelihood function.

553 Consider model parameterization  $z$  and data  $x$  generated from some theoretical model simulator  
 554 represented as  $p(x | z)$ , which may be deterministic or stochastic. Theoretical models usually have  
 555 known sampling procedures for simulating activity given a circuit parameterization, yet often lack  
 556 an explicit likelihood function due to the nonlinearities and dynamics. With EPI, a distribution  
 557 on parameters  $z$  is learned, that yields an emergent property of interest  $\mathcal{B}$ ,

$$\mathcal{B} \leftrightarrow E_{z \sim q_\theta} [E_{x \sim p(x|z)} [T(x)]] = \mu \quad (12)$$

558 by making an approximation  $q_\theta(z)$  to  $p(z | \mathcal{B})$  (see Section A.1.5). So, over the DSN distribution  
 559  $q_\theta(z)$  of model  $p(x | z)$  for behavior  $\mathcal{B}$ , the emergent properties  $T(x)$  are constrained in expectation  
 560 to  $\mu$ .

561 In deep probability distributions, a simple random variable  $w \sim q_0$  is mapped deterministically via  
 562 a function  $f_\theta$  parameterized by a neural network to the support of the distribution of interest where  
 563  $z = f_\theta(w) = f_l(\dots f_1(w))$ . Given a theoretical model  $p(x | z)$  and some behavior of interest  $\mathcal{B}$ , the  
 564 deep probability distributions are trained by optimizing the neural network parameters  $\theta$  to find a  
 565 good approximation  $q_\theta^*$  within the deep variational family  $Q$  to  $p(z | \mathcal{B})$ .

566 In most settings (especially those relevant to theoretical neuroscience) the likelihood of the behavior  
 567 with respect to the model parameters  $p(T(x) | z)$  is unknown or intractable, requiring an alternative  
 568 to stochastic gradient variational Bayes [5] or black box variational inference[53]. These types  
 569 of methods called likelihood-free variational inference (LFVI, [20]) skate around the intractable  
 570 likelihood function in situations where there is a differentiable simulator. Akin to LFVI, DSNs are  
 571 optimized with the following objective for a given theoretical model, emergent property statistics  
 572  $T(x)$ , and emergent property constraints  $\mu$ :

$$\begin{aligned} q_\theta^*(z) &= \operatorname{argmax}_{q_\theta \in Q} H(q_\theta(z)) \\ \text{s.t. } E_{z \sim q_\theta} [E_{x \sim p(x|z)} [T(x)]] &= \mu \end{aligned} \quad (13)$$

573 Optimizing this objective is a technological accomplishment in its own right, the details of which  
 574 we elaborate in Section A.1.2. Before going through those details, we ground this optimization in  
 575 a toy example.

576 **A.1.1 Example: 2D LDS**

577 To gain intuition for EPI, consider two-dimensional linear dynamical systems,  $\tau \dot{x} = Ax$  with

$$A = \begin{bmatrix} a_1 & a_2 \\ a_3 & a_4 \end{bmatrix}$$

578 that produce a band of oscillations. To do EPI with the dynamics matrix elements as the free  
 579 parameters  $z = [a_1, a_2, a_3, a_4]$ , and fixing  $\tau = 1$ , such that the posterior yields a band of oscillations,  
 580 the emergent property statistics  $T(x)$  are chosen to contain the first- and second-moments of the  
 581 oscillatory frequency  $\omega$  and the growth/decay factor  $d$  of the oscillating system. To learn the  
 582 distribution of real entries of  $A$  that yield a distribution of  $d$  with mean zero with variance  $0.25^2$ ,  
 583 and oscillation frequency  $\omega$  with mean 1 Hz with variance  $(0.1\text{Hz})^2$ , then we would select the real  
 584 part of the complex conjugate eigenvalues  $\text{real}(\lambda_1) = d$  (via an arbitrary choice of eigenvalue of the  
 585 dynamics matrix  $\lambda_1$ ) and the positive imaginary component of one of the eigenvalues  $\text{imag}(\lambda_1) =$   
 586  $2\pi\omega$  as the emergent property statistics. Those emergent property statistics are then constrained  
 587 to

$$\mu = E \begin{bmatrix} \text{real}(\lambda_1) \\ \text{imag}(\lambda_1) \\ (\text{real}(\lambda_1) - 0)^2 \\ (\text{imag}(\lambda_1) - 2\pi\omega)^2 \end{bmatrix} = \begin{bmatrix} 0.0 \\ 2\pi\omega \\ 0.25^2 \\ (2\pi 0.1)^2 \end{bmatrix} \quad (14)$$

588 where  $\omega = 1\text{Hz}$ . Unlike the models we study in the paper which calculate  $E_{x \sim p(x|z)} [T(x)]$  via  
 589 forward simulation, we have a closed form for the eigenvalues of the dynamics matrix.  $\lambda$  can be  
 590 calculated using the quadratic formula:

$$\lambda = \frac{\left(\frac{a_1+a_4}{\tau}\right) \pm \sqrt{\left(\frac{a_1+a_4}{\tau}\right)^2 + 4\left(\frac{a_2a_3-a_1a_4}{\tau}\right)}}{2} \quad (15)$$

591 where  $\lambda_1$  is the eigenvalue of  $\frac{1}{\tau}A$  with greatest real part. Even though  $E_{x \sim p(x|z)} [T(x)]$  is calculable  
 592 directly via a closed form function and does not require simulation, we cannot derive the distribution  
 593  $q_\theta^*$  directly. This is due to the formally hard problem of the backward mapping: finding the natural  
 594 parameters  $\eta$  from the mean parameters  $\mu$  of an exponential family distribution [54]. Instead, we  
 595 can use EPI to learn the linear system parameters producing such a band of oscillations (Fig. S2B).

596 Even this relatively simple system has nontrivial (though intuitively sensible) structure in the  
 597 parameter distribution. To validate our method (further than that of the underlying technology  
 598 on a ground truth solution [19]) we can analytically derive the contours of the probability density

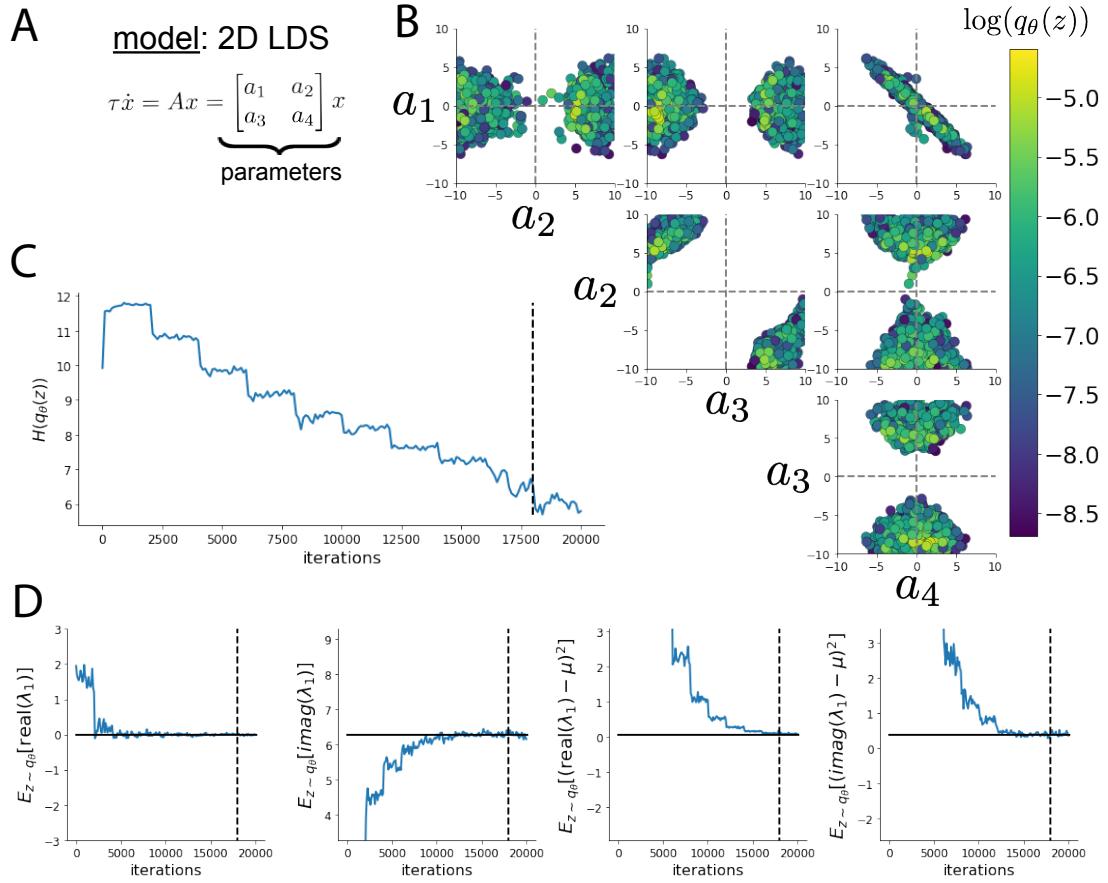


Fig. S2: A. Two-dimensional linear dynamical system model, where real entries of the dynamics matrix  $A$  are the parameters. B. The DSN distribution for a 2D LDS with  $\tau = 1$  that produces an average of 1Hz oscillations with some small amount of variance. C. Entropy throughout the optimization. At the beginning of each augmented Lagrangian epoch (5,000 iterations), the entropy dips 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 move closer to their constraints.

599 from the emergent property statistics and values (Fig. S3). In the  $a_1 - a_4$  plane, is a black line  
 600 at  $\text{real}(\lambda_1) = \frac{a_1 + a_4}{2} = 0$ , a dotted black line at the standard deviation  $\text{real}(\lambda_1) = \frac{a_1 + a_4}{2} \pm 1$ , and a  
 601 grey line at twice the standard deviation  $\text{real}(\lambda_1) = \frac{a_1 + a_4}{2} \pm 2$  (Fig. S3A). Here the lines denote the  
 602 set of solutions at fixed behaviors, which overlay the posterior obtained through EPI. The learned  
 603 DSN distribution precisely reflects the desired statistical constraints and model degeneracy in the  
 604 sum of  $a_1$  and  $a_4$ . Intuitively, the parameters equivalent with respect to emergent property statistic  
 605  $\text{real}(\lambda_1)$  have similar log densities.

606 To explain the structure in the bimodality of the DSN posterior, we can look at the imaginary  
 607 component of  $\lambda_1$ . When  $\text{real}(\lambda_1) = \frac{a_1 + a_4}{2} = 0$ , we have

$$\text{imag}(\lambda_1) = \begin{cases} \sqrt{\frac{a_1 a_4 - a_2 a_3}{\tau}}, & \text{if } a_1 a_4 < a_2 a_3 \\ 0 & \text{otherwise} \end{cases} \quad (16)$$

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

$$\text{imag}(\lambda_1)^2 = a_1 a_4 - a_2 a_3 \quad (17)$$

610 Since we constrained  $E_{q_\theta}[\text{imag}(\lambda)] = 2\pi$  (with  $\omega = 1$ ), we can plot contours of the equation  
 611  $\text{imag}(\lambda_1)^2 = a_1 a_4 - a_2 a_3 = (2\pi)^2$  for various  $a_1 a_4$  (Fig. S3A). If  $\sigma_{1,4} = E_{q_\theta}(|a_1 a_4 - E_{q_\theta}[a_1 a_4]|)$ ,  
 612 then we plot the contours as  $a_1 a_4 = 0$  (black),  $a_1 a_4 = -\sigma_{1,4}$  (black dotted), and  $a_1 a_4 = -2\sigma_{1,4}$   
 613 (grey dotted) (Fig. S3B). This validates the curved structure of the inferred distribution learned  
 614 through EPI. We take steps in negative standard deviation of  $a_1 a_4$  (dotted and gray lines), since  
 615 there are few positive values  $a_1 a_4$  in the posterior. Subtler model-behavior combinations will have  
 616 even more complexity, further motivating the use of EPI for understanding these systems. Indeed,  
 617 we sample a distribution of systems oscillating near 1Hz (Fig. S4).

### 618 A.1.2 Augmented Lagrangian optimization

619 To optimize  $q_\theta(z)$  in Equation 13, the constrained optimization is performed using the augmented  
 620 Lagrangian method. The following objective is minimized:

$$L(\theta; \alpha, c) = -H(q_\theta) + \alpha^\top \delta(\theta) + \frac{c}{2} \|\delta(\theta)\|^2 \quad (18)$$

621 where  $\delta(\theta) = E_{z \sim q_\theta} [E_{x \sim p(x|z)} [T(x) - \mu]]$ ,  $\alpha \in \mathcal{R}^m$  are the Lagrange multipliers and  $c$  is the penalty  
 622 coefficient. For a fixed  $(\alpha, c)$ ,  $\theta$  is optimized with stochastic gradient descent. A low value of  $c$  is

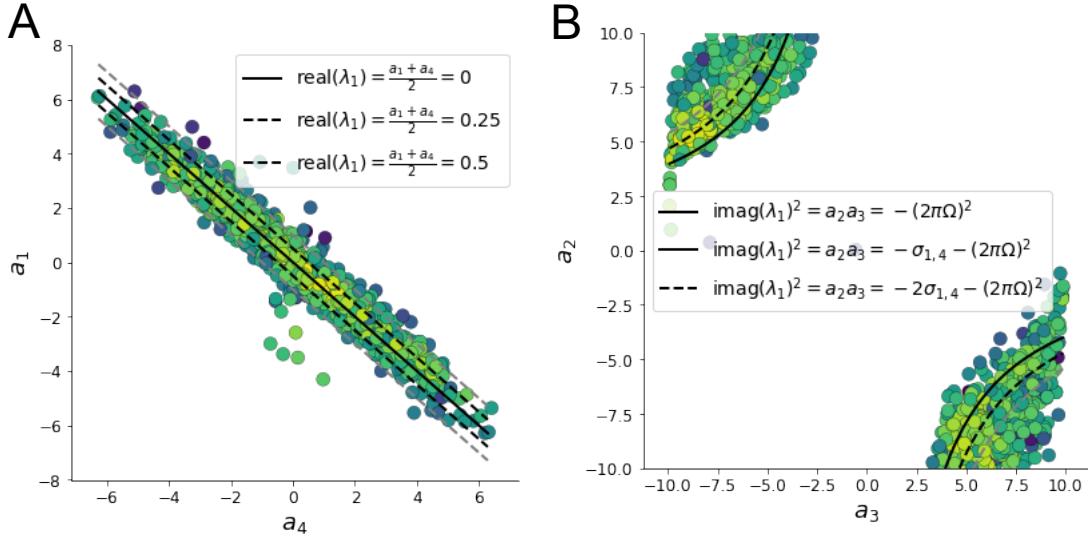


Fig. S3: A. Probability contours in the  $a_1 - a_4$  plane can be derived from the relationship to emergent property statistic of growth/decay factor. B. Probability contours in the  $a_2 - a_3$  plane can be derived from relationship to the emergent property statistic of oscillation frequency.

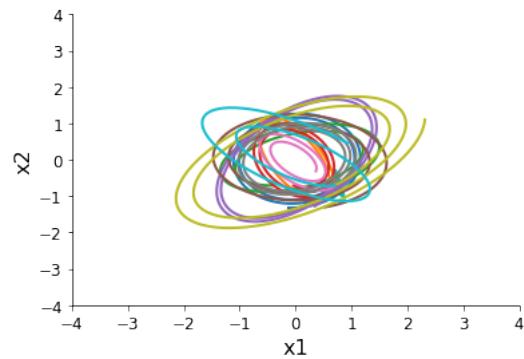


Fig. S4: Sampled dynamical system trajectories from the EPI distribution. Each trajectory is initialized at  $x(0) = \left[ \frac{\sqrt{2}}{2} \quad -\frac{\sqrt{2}}{2} \right]$ .

623 used initially, and increased during each augmented Lagrangian epoch – a period of optimization  
 624 with fixed  $\alpha$  and  $c$  for a given number of stochastic optimization iterations. Similarly,  $\alpha$  is tuned  
 625 each epoch based on the constraint violations. For the linear 2-dimensional system (Fig. S2C)  
 626 optimization hyperparameters are initialized to  $c_1 = 10^{-4}$  and  $\alpha_1 = 0$ . The penalty coefficient  
 627 is updated based on a hypothesis test regarding the reduction in constraint violation. The p-  
 628 value of  $E[\|\delta(\theta_{k+1})\|] > \gamma E[\|\delta(\theta_k)\|]$  is computed, and  $c_{k+1}$  is updated to  $\beta c_k$  with probability  
 629  $1 - p$ . Throughout the project,  $\beta = 4.0$  and  $\gamma = 0.25$  is used. The other update rule is  $\alpha_{k+1} =$   
 630  $\alpha_k + c_k \frac{1}{n} \sum_{i=1}^n (T(x^{(i)}) - \mu)$ . In this example, each augmented Lagrangian epoch ran for 2,000  
 631 iterations. We consider the optimization to have converged when a null hypothesis test of constraint  
 632 violations being zero is accepted for all constraints at a significance threshold 0.05. This is the dotted  
 633 line on the plots below depicting the optimization cutoff of EPI optimization for the 2-dimensional  
 634 linear system. If the optimization is left to continue running, entropy usually decreases, and  
 635 structural pathologies in the distribution may be introduced.

636 The intention is that  $c$  and  $\alpha$  start at values encouraging entropic growth early in optimization.  
 637 Then, as they increase in magnitude with each training epoch, the constraint satisfaction terms are  
 638 increasingly weighted, resulting in a decrease in entropy. Rather than using a naive initialization,  
 639 before EPI, we optimize the deep probability distribution parameters to generate samples of an  
 640 isotropic Gaussian of a selected variance, such as 1.0 for the 2D LDS example. This provides a  
 641 convenient starting point, whose level of entropy is controlled by the user.

### 642 A.1.3 Normalizing flows

643 Since we are optimizing parameters  $\theta$  of our deep probability distribution with respect to the  
 644 entropy, we will need to take gradients with respect to the log-density of samples from the deep  
 645 probability distribution.

$$H(q_\theta(z)) = \int -q_\theta(z) \log(q_\theta(z)) dz = E_{z \sim q_\theta} [-\log(q_\theta(z))] = E_{w \sim q_0} [-\log(q_\theta(f_\theta(w)))] \quad (19)$$

$$\nabla_\theta H(q_\theta(z)) = E_{w \sim q_0} [-\nabla_\theta \log(q_\theta(f_\theta(w)))] \quad (20)$$

647 Deep probability models typically consist of several layers of fully connected neural networks.  
 648 When each neural network layer is restricted to be a bijective function, the sample density can be  
 649 calculated using the change of variables formula at each layer of the network. For  $z' = f(z)$ ,

$$q(z') = q(f^{-1}(z')) \left| \det \frac{\partial f^{-1}(z')}{\partial z'} \right| = q(z) \left| \det \frac{\partial f(z)}{\partial z} \right|^{-1} \quad (21)$$

650 However, this computation has cubic complexity in dimensionality for fully connected layers. By  
 651 restricting our layers to normalizing flows [16] – bijective functions with fast log determinant ja-  
 652 cobian computations, we can tractably optimize deep generative models with objectives that are a  
 653 function of sample density, like entropy. Most of our analyses use real NVP [55], which have proven  
 654 effective in our architecture searches, and have the advantageous features of fast sampling and fast  
 655 density evaluation.

656 **A.1.4 Related work**

657 (To come)

658

659 **A.1.5 Emergent property inference as variational inference in an exponential family**

660 (To come)

661

662 **A.2 Theoretical models**

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

665 **A.2.1 Stomatogastric ganglion**

666 Each neuron's membrane potential  $x_m(t)$  is the solution of the following differential equation.

$$C_m \frac{dx_m}{dt} = -[h_{leak}(x; z) + h_{Ca}(x; z) + h_K(x; z) + h_{hyp}(x; z) + h_{elec}(x; z) + h_{syn}(x; z)] \quad (22)$$

667 The membrane potential of each neuron is affected by the leak, calcium, potassium, hyperpolariza-  
 668 tion, electrical and synaptic currents, respectively. The capacitance of the cell membrane was set to  
 669  $C_m = 1nF$ . Each current is a function of the neuron's membrane potential  $x_m$  and the parameters  
 670 of the circuit such as  $g_{el}$  and  $g_{syn}$ , whose effect on the circuit is considered in the motivational

example of EPI in Fig. 1. Specifically, the currents are the difference in the neuron's membrane potential and that current type's reversal potential multiplied by a conductance:

$$h_{leak}(x; z) = g_{leak}(x_m - V_{leak}) \quad (23)$$

673

$$h_{elec}(x; z) = g_{el}(x_m^{post} - x_m^{pre}) \quad (24)$$

674

$$h_{syn}(x; z) = g_{syn}S_\infty^{pre}(x_m^{post} - V_{syn}) \quad (25)$$

675

$$h_{Ca}(x; z) = g_{Ca}M_\infty(x_m - V_{Ca}) \quad (26)$$

676

$$h_K(x; z) = g_KN(x_m - V_K) \quad (27)$$

677

$$h_{hyp}(x; z) = g_hH(x_m - V_{hyp}) \quad (28)$$

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

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

$$M_\infty = 0.5 \left( 1 + \tanh \left( \frac{x_m - v_1}{v_2} \right) \right) \quad (29)$$

685

$$\frac{dN}{dt} = \lambda_N(N_\infty - N) \quad (30)$$

686

$$N_\infty = 0.5 \left( 1 + \tanh \left( \frac{x_m - v_3}{v_4} \right) \right) \quad (31)$$

687

$$\lambda_N = \phi_N \cosh \left( \frac{x_m - v_3}{2v_4} \right) \quad (32)$$

688

$$\frac{dH}{dt} = \frac{(H_\infty - H)}{\tau_h} \quad (33)$$

689

$$H_\infty = \frac{1}{1 + \exp \left( \frac{x_m + v_5}{v_6} \right)} \quad (34)$$

690

$$\tau_h = 272 - \left( \frac{-1499}{1 + \exp \left( \frac{-x_m + v_7}{v_8} \right)} \right) \quad (35)$$

where we set  $v_1 = 0mV$ ,  $v_2 = 20mV$ ,  $v_3 = 0mV$ ,  $v_4 = 15mV$ ,  $v_5 = 78.3mV$ ,  $v_6 = 10.5mV$ ,  $v_7 = -42.2mV$ ,  $v_8 = 87.3mV$ ,  $v_9 = 5mV$ , and  $v_{th} = -25mV$ . These are the same parameter values used in [22].

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

$$S_\infty = \frac{1}{1 + \exp\left(\frac{v_{th} - x_m}{v_0}\right)} \quad (36)$$

695 When the dynamic gating variables are considered, this is actually a 15-dimensional nonlinear  
696 dynamical system.

697 In order to measure the frequency of the hub neuron during EPI, the STG model was simulated  
698 for  $T = 500$  time steps of  $dt = 25ms$ . In EPI, since gradients are taken through the simulation  
699 process, the number of time steps are kept as modest if possible. The chosen  $dt$  and  $T$  were the  
700 most computationally convenient choices yielding accurate frequency measurement.

701 Our original approach to measuring frequency was to take the max of the fast Fourier transform  
702 (FFT) of the simulated time series. There are a few key considerations here. One is resolution  
703 in frequency space. Each FFT entry will correspond to a signal frequency of  $\frac{F_s k}{N}$ , where  $N$  is  
704 the number of samples used for the FFT,  $F_s = \frac{1}{dt}$ , and  $k \in [0, 1, \dots, N - 1]$ . Our resolution is  
705 improved by increasing  $N$  and decreasing  $dt$ . Increasing  $N = T - b$ , where  $b$  is some fixed number  
706 of buffer burn-in initialization samples, necessitates an increase in simulation time steps  $T$ , which  
707 directly increases computational cost. Increasing  $F_s$  (decreasing  $dt$ ) increases system approximation  
708 accuracy, but requires more time steps before a full cycle is observed. At the level of  $dt = 0.025$ ,  
709 thousands of temporal samples were required for resolution of .01Hz. These challenges in frequency  
710 resolution with the discrete Fourier transform motivated the use of an alternative basis of complex  
711 exponentials. Instead, we used a basis of complex exponentials with frequencies from 0.0-1.0 Hz at  
712 0.01Hz resolution,  $\Phi = [0.0, 0.01, \dots, 1.0]^\top$

713 Another consideration was that the frequency spectra of the hub neuron has several peaks. This  
714 was due to high-frequency sub-threshold activity. The maximum frequency was often not the firing  
715 frequency. Accordingly, subthreshold activity was set to zero, and the whole signal was low-pass  
716 filtered with a moving average window of length 20. The signal was subsequently mean centered.  
717 After this pre-processing, the maximum frequency in the filter bank accurately reflected the firing  
718 frequency.

719 Finally, to differentiate through the maximum frequency identification step, we used a sum-of-  
720 powers normalization strategy: Let  $\mathcal{X}_i \in \mathcal{C}^{|\Phi|}$  be the complex exponential filter bank dot products  
721 with the signal  $x_i \in \mathcal{R}^N$ , where  $i \in \{\text{f1}, \text{f2}, \text{hub}, \text{s1}, \text{s2}\}$ . The “frequency identification” vector is

$$u_i = \frac{|\mathcal{X}_i|^\alpha}{\sum_{k=1}^N |\mathcal{X}_i(k)|^\alpha} \quad (37)$$

722 The frequency is then calculated as  $\omega = u_i^\top \Phi$  with  $\alpha = 100$ .

723 Network syncing, like all other emergent properties in this work, are defined by the emergent  
 724 property statistics and values. The emergent property statistics are the first- and second-moments  
 725 of the firing frequencies. The first moments are set to 0.542Hz, while the second moments are set  
 726 to 0.025Hz<sup>2</sup>.

$$E \begin{bmatrix} \omega_{f1} \\ \omega_{f2} \\ \omega_{hub} \\ \omega_{s1} \\ \omega_{s2} \\ (\omega_{f1} - 0.542)^2 \\ (\omega_{f2} - 0.542)^2 \\ (\omega_{hub} - 0.542)^2 \\ (\omega_{s1} - 0.542)^2 \\ (\omega_{s2} - 0.542)^2 \end{bmatrix} = \begin{bmatrix} 0.542 \\ 0.542 \\ 0.542 \\ 0.542 \\ 0.542 \\ 0.025^2 \\ 0.025^2 \\ 0.025^2 \\ 0.025^2 \\ 0.025^2 \end{bmatrix} \quad (38)$$

727 For EPI in Fig 2C, we used a real NVP architecture with two coupling layers. Each coupling layer  
 728 had two hidden layers of 10 units each, and we mapped onto a support of  $z \in \left[ \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 10 \\ 8 \end{bmatrix} \right]$ . We  
 729 have shown the EPI optimization that converged with maximum entropy across 2 random seeds  
 730 and augmented Lagrangian coefficient initializations of  $c_0=0$ , 2, and 5.

731 **A.2.2 Primary visual cortex**

732 The dynamics of each neural populations average rate  $x = \begin{bmatrix} x_E \\ x_P \\ x_S \\ x_V \end{bmatrix}$  are given by:

$$\tau \frac{dx}{dt} = -x + [Wx + h]_+^n \quad (39)$$

733 Some neuron-types largely lack synaptic projections to other neuron-types [56], and it is popular

<sup>734</sup> to only consider a subset of the effective connectivities [23].

$$W = \begin{bmatrix} W_{EE} & W_{EP} & W_{ES} & 0 \\ W_{PE} & W_{PP} & W_{PS} & 0 \\ W_{SE} & 0 & 0 & W_{SV} \\ W_{VE} & W_{VP} & W_{VS} & 0 \end{bmatrix} \quad (40)$$

<sup>735</sup> By consolidating information from many experimental datasets, Billeh et al. [40] produce estimates  
<sup>736</sup> of the synaptic strength (in mV)

$$M = \begin{bmatrix} 0.36 & 0.48 & 0.31 & 0.28 \\ 1.49 & 0.68 & 0.50 & 0.18 \\ 0.86 & 0.42 & 0.15 & 0.32 \\ 1.31 & 0.41 & 0.52 & 0.37 \end{bmatrix} \quad (41)$$

<sup>737</sup> and connection probability

$$C = \begin{bmatrix} 0.16 & 0.411 & 0.424 & 0.087 \\ 0.395 & .451 & 0.857 & 0.02 \\ 0.182 & 0.03 & 0.082 & 0.625 \\ 0.105 & 0.22 & 0.77 & 0.028 \end{bmatrix} \quad (42)$$

<sup>738</sup> Multiplying these connection probabilities and synaptic efficacies gives us an effective connectivity  
<sup>739</sup> matrix:

$$W_{\text{full}} = C \odot M = \begin{bmatrix} 0.16 & 0.411 & 0.424 & 0.087 \\ 0.395 & .451 & 0.857 & 0.02 \\ 0.182 & 0.03 & 0.082 & 0.625 \\ 0.105 & 0.22 & 0.77 & 0.028 \end{bmatrix} \quad (43)$$

<sup>740</sup> From use the entries of this full effective connectivity matrix that are not considered to be ineffectual.  
<sup>741</sup>

<sup>742</sup> We look at how this four-dimensional nonlinear dynamical model of V1 responds to different inputs,  
<sup>743</sup> and compare the predictions of the linear response to the approximate posteriors obtained through  
<sup>744</sup> EPI. The input to the system is the sum of a baseline input  $b = [1 \ 1 \ 1 \ 1]^T$  and a differential  
<sup>745</sup> input  $dh$ :

$$h = b + dh \quad (44)$$

<sup>746</sup> All simulations of this system had  $T = 100$  time points, a time step  $dt = 5\text{ms}$ , and time constant  
<sup>747</sup>  $\tau = 20\text{ms}$ . And the system was initialized to a random draw  $x(0)_i \sim \mathcal{N}(1, 0.01)$ .

<sup>748</sup> We can describe the dynamics of this system more generally by

$$\dot{x}_i = -x_i + f(u_i) \quad (45)$$

<sup>749</sup> where the input to each neuron is

$$u_i = \sum_j W_{ij} x_j + h_i \quad (46)$$

<sup>750</sup> Let  $F_{ij} = \gamma_i \delta(i, j)$ , where  $\gamma_i = f'(u_i)$ . Then, the linear response is

$$\frac{dx_{ss}}{dh} = F(W \frac{dx_{ss}}{dh} + I) \quad (47)$$

<sup>751</sup> which is calculable by

$$\frac{dx_{ss}}{dh} = (F^{-1} - W)^{-1} \quad (48)$$

<sup>752</sup> The emergent property we considered was the first and second moments of the change in rate  $dx$   
<sup>753</sup> between the baseline input  $h = b$  and  $h = b + dh$ . We use the following notation to indicate that  
<sup>754</sup> the emergent property statistics were set to the following values:

$$\mathcal{B}(\alpha, y) \leftrightarrow E \begin{bmatrix} dx_{\alpha,ss} \\ (dx_{\alpha,ss} - y)^2 \end{bmatrix} = \begin{bmatrix} y \\ 0.01^2 \end{bmatrix} \quad (49)$$

<sup>755</sup> In the final analysis for this model, we sweep the input one neuron at a time away from the mode  
<sup>756</sup> of each inferred distributions  $dh^* = z^* = \text{argmax}_z \log q_\theta(z \mid \mathcal{B}(\alpha, 0.1))$ . The differential responses  
<sup>757</sup>  $dx_{\alpha,ss}$  are examined at perturbed inputs  $h = b + dh^* + \Delta h_\alpha u_\alpha$  where  $u_\alpha$  is a unit vector in the  
<sup>758</sup> dimension of  $\alpha$  and  $\Delta h_\alpha \in [-15, 15]$ .

<sup>759</sup> For each  $\mathcal{B}(\alpha, y)$  with  $\alpha \in \{E, P, S, V\}$  and  $y \in \{0.1, 0.5\}$ , we ran EPI with five different random  
<sup>760</sup> initial seeds using an architecture of four coupling layers, each with two hidden layers of 10 units.  
<sup>761</sup> We set  $c_0 = 10^5$ . The support of the learned distribution was restricted to  $z_i \in [-5, 5]$ .

### <sup>762</sup> A.2.3 Superior colliculus

<sup>763</sup> There are four total units: two in each hemisphere corresponding to the Pro/Contra and Anti/Ipsi  
<sup>764</sup> populations. Each unit has an activity ( $x_i$ ) and internal variable ( $u_i$ ) related by

$$x_i(t) = \left( \frac{1}{2} \tanh \left( \frac{v_i(t) - \epsilon}{\zeta} \right) + \frac{1}{2} \right) \quad (50)$$

<sup>765</sup>  $\epsilon = 0.05$  and  $\zeta = 0.5$  control the position and shape of the nonlinearity, respectively.

<sup>766</sup> We can order the elements of  $x_i$  and  $v_i$  into vectors  $x$  and  $v$  with elements

$$x = \begin{bmatrix} x_{LP} \\ x_{LA} \\ x_{RP} \\ x_{RA} \end{bmatrix} \quad v = \begin{bmatrix} v_{LP} \\ v_{LA} \\ v_{RP} \\ v_{RA} \end{bmatrix} \quad (51)$$

<sup>767</sup> The internal variables follow dynamics:

$$\tau \frac{dv}{dt} = -v + Wx + h + \sigma dB \quad (52)$$

<sup>768</sup> with time constant  $\tau = 0.09s$  and Gaussian noise  $\sigma dB$  controlled by the magnitude of  $\sigma = 1.0$ . The  
<sup>769</sup> weight matrix has 8 parameters  $sW_P$ ,  $sW_A$ ,  $vW_{PA}$ ,  $vW_{AP}$ ,  $hW_P$ ,  $hW_A$ ,  $dW_{PA}$ , and  $dW_{AP}$  (Fig.  
<sup>770</sup> 4B).

$$W = \begin{bmatrix} sW_P & vW_{PA} & hW_P & dW_{PA} \\ vW_{AP} & sW_A & dW_{AP} & hW_A \\ hW_P & dW_{PA} & sW_P & vW_{PA} \\ dW_{AP} & hW_A & vW_{AP} & sW_A \end{bmatrix} \quad (53)$$

<sup>771</sup> The system receives five inputs throughout each trial, which has a total length of 1.8s.

$$h = h_{\text{rule}} + h_{\text{choice-period}} + h_{\text{light}} \quad (54)$$

<sup>772</sup> There are rule-based inputs depending on the condition,

$$h_{P,\text{rule}}(t) = \begin{cases} I_{P,\text{rule}} \begin{bmatrix} 1 & 0 & 0 & 1 \end{bmatrix}^\top, & \text{if } t \leq 1.2s \\ 0, & \text{otherwise} \end{cases} \quad (55)$$

$$h_{A,\text{rule}}(t) = \begin{cases} I_{A,\text{rule}} \begin{bmatrix} 0 & 1 & 1 & 0 \end{bmatrix}^\top, & \text{if } t \leq 1.2s \\ 0, & \text{otherwise} \end{cases} \quad (56)$$

<sup>774</sup> a choice-period input,

$$h_{\text{choice}}(t) = \begin{cases} I_{\text{choice}} \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix}^\top, & \text{if } t > 1.2s \\ 0, & \text{otherwise} \end{cases} \quad (57)$$

<sup>775</sup> and an input to the right or left-side depending on where the light stimulus is delivered.

$$h_{\text{light}}(t) = \begin{cases} I_{\text{light}} \begin{bmatrix} 1 & 1 & 0 & 0 \end{bmatrix}^\top, & \text{if } t > 1.2s \text{ and Left} \\ I_{\text{light}} \begin{bmatrix} 0 & 0 & 1 & 1 \end{bmatrix}^\top, & \text{if } t > 1.2s \text{ and Right} \\ 0, & t \leq 1.2s \end{cases} \quad (58)$$

776 The input parameterization was fixed to  $I_{P,\text{rule}} = 10$ ,  $I_{A,\text{rule}} = 10$ ,  $I_{\text{choice}} = 2$ , and  $I_{\text{light}} = 1$   
 777 To produce a Bernoulli rate of  $p_{LP}$  in the Left, Pro condition (we can generalize this to either cue,  
 778 or stimulus condition), let  $\hat{p}_i$  be the empirical average steady state (ss) response (final  $x_{LP}$  at end  
 779 of task) over  $M=500$  Gaussian noise draws for a given SC model parameterization  $z_i$ :

$$\hat{p}_i = E_{\sigma dB} [x_{LP,ss} | s = L, c = P, z_i] = \frac{1}{M} \sum_{j=1}^M x_{LP,ss}(s = L, c = P, z_i, \sigma dB_j) \quad (59)$$

780 For the first constraint, the average over posterior samples (from  $q_\theta(z)$ ) to be  $p_{LP}$ :

$$E_{z_i \sim q_\phi} [E_{\sigma dB} [x_{LP,ss} | s = L, c = P, z_i]] = E_{z_i \sim q_\phi} [\hat{p}_i] = p_{LP} \quad (60)$$

781 We can then ask that the variance of the steady state responses across Gaussian draws, is the  
 782 Bernoulli variance for the empirical rate  $\hat{p}_i$ .

$$E_{z \sim q_\phi} [\sigma_{err}^2] = 0 \quad (61)$$

783

$$\sigma_{err}^2 = Var_{\sigma dB} [x_{LP,ss} | s = L, c = P, z_i] - \hat{p}_i(1 - \hat{p}_i) \quad (62)$$

784 We have an additional constraint that the Pro neuron on the opposite hemisphere should have the  
 785 opposite value. We can enforce this with a final constraint:

$$E_{z \sim q_\phi} [d_P] = 1 \quad (63)$$

786

$$E_{\sigma dB} [(x_{LP,ss} - x_{RP,ss})^2 | s = L, c = P, z_i] \quad (64)$$

787 We refer to networks obeying these constraints as Bernoulli, winner-take-all networks. Since the  
 788 maximum variance of a random variable bounded from 0 to 1 is the Bernoulli variance ( $\hat{p}(1 - \hat{p})$ ),  
 789 and the maximum squared difference between two variables bounded from 0 to 1 is 1, we do not  
 790 need to control the second moment of these test statistics. In reality, these variables are dynamical  
 791 system states and can only exponentially decay (or saturate) to 0 (or 1), so the Bernoulli variance  
 792 error and squared difference constraints can only be undershot. This is important to be mindful  
 793 of when evaluating the convergence criteria. Instead of using our usual hypothesis testing criteria  
 794 for convergence to the emergent property, we set a slack variable threshold for these technically  
 795 infeasible constraints to 0.05.

796 Training DSNs to learn distributions of dynamical system parameterizations that produce Bernoulli  
 797 responses at a given rate (with small variance around that rate) was harder to do than expected.

798 There is a pathology in this optimization setup, where the learned distribution of weights is bimodal  
 799 attributing a fraction  $p$  of the samples to an expansive mode (which always sends  $x_{LP}$  to 1), and a  
 800 fraction  $1 - p$  to a decaying mode (which always sends  $x_{LP}$  to 0). This pathology was avoided using  
 801 an inequality constraint prohibiting parameter samples that resulted in low variance of responses  
 802 across noise.

803 In total, the emergent property of rapid task switching accuracy at level  $p$  was defined as

$$\mathcal{B}(p) \leftrightarrow \begin{bmatrix} \hat{p}_P \\ \hat{p}_A \\ (\hat{p}_P - p)^2 \\ (\hat{p}_A - p)^2 \\ \sigma_{P,err}^2 \\ \sigma_{A,err}^2 \\ d_P \\ d_A \end{bmatrix} = \begin{bmatrix} p \\ p \\ 0.15^2 \\ 0.15^2 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} \quad (65)$$

804 For each accuracy level  $p$ , we ran EPI for 10 different random seeds and selected the maximum  
 805 entropy solution using an architecture of 10 planar flows with  $c_0 = 2$ . The support of  $z$  was  $\mathcal{R}^8$ .

#### 806 A.2.4 Rank-1 RNN

807 Recent work establishes a link between RNN connectivity weights and the resulting dynamical  
 808 responses of the network, using dynamic mean field theory (DMFT) [25]. Specifically, DMFT  
 809 describes the properties of activity in infinite-size neural networks given a distribution on the  
 810 connectivity weights. In such a model, the connectivity of a rank-1 RNN (which was sufficient for  
 811 our task), has weight matrix  $W$ , which is the sum of a random component with strength determined  
 812 by  $g$  and a structured component determined by the outer product of vectors  $m$  and  $n$ :

$$W = g\chi + \frac{1}{N}mn^\top, \quad (66)$$

813 where the activity  $x$  evolves as and  $I(t)$  is some input,  $\phi$  is the tanh nonlinearity, and  $\chi_{ij} \sim \mathcal{N}(0, \frac{1}{N})$ .  
 814 The entries of  $m$  and  $n$  are drawn from Gaussian distributions  $m_i \sim \mathcal{N}(M_m, 1)$  and  $n_i \sim \mathcal{N}(M_n, 1)$ .  
 815 From such a parameterization, this theory produces consistency equations for the dynamic mean  
 816 field variables in terms of parameters like  $g$ ,  $M_m$ , and  $M_n$ , which we study in Section 3.5. That  
 817 is the dynamic mean field variables (e.g. the activity along a vector  $\kappa_v$ , the total variance

818  $\Delta_0$ , structured variance  $\Delta_\infty$ , and the chaotic variance  $\Delta_T$ ) are written as functions of one another  
819 in terms of connectivity parameters. The values of these variables can be used obtained using a  
820 nonlinear system of equations solver. These dynamic mean field variables are then cast as task-  
821 relevant variables with respect to the context of the provided inputs. Mastrogiuseppe et al. designed  
822 low-rank RNN connectivities via minimalist connectivity parameters to solve canonical tasks from  
823 behavioral neuroscience.

824 We consider the DMFT equation solver as a black box that takes in a low-rank parameterization  $z$   
825 (e.g.  $z = [g \quad M_m \quad M_n]$ ) and outputs the values of the dynamic mean field variables, of which we  
826 cast  $\kappa_w$  and  $\Delta_T$  as task-relevant variables  $\mu_{\text{post}}$  and  $\sigma_{\text{post}}^2$  in the Gaussian posterior conditioning  
827 toy example. Importantly, the solution produced by the solver is differentiable with respect to the  
828 input parameters, allowing us to use DMFT to calculate the emergent property statistics in EPI  
829 to learn distributions on such connectivity parameters of RNNs that execute tasks.

830 Specifically, we solve for the mean field variables  $\kappa_w$ ,  $\kappa_n$ ,  $\Delta_0$  and  $\Delta_\infty$ , where the readout is nominally  
831 chosen to point in the unit orthant  $w = [1 \quad \dots \quad 1]^\top$ . The consistency equations for these variables  
832 in the presence of an constant input  $I(t) = y - (n - M_n)$  can be derived following [25] are

$$\begin{aligned} \kappa_w &= F(\kappa_w, \kappa_n, \Delta_0, \Delta_\infty) = M_m \kappa_n + y \\ \kappa_n &= G(\kappa_w, \kappa_n, \Delta_0, \Delta_\infty) = M_n \langle [\phi_i] \rangle + \langle [\phi'_i] \rangle \\ \frac{\Delta_0^2 - \Delta_\infty^2}{2} &= H(\kappa_w, \kappa_n, \Delta_0, \Delta_\infty) = g^2 \left( \int \mathcal{D}z \Phi^2(\kappa_w + \sqrt{\Delta_0} z) - \int \mathcal{D}z \int \mathcal{D}x \Phi(\kappa_w + \sqrt{\Delta_0 - \Delta_\infty} x + \sqrt{\Delta_\infty} z) \right) \\ &\quad + (\kappa_n^2 + 1)(\Delta_0 - \Delta_\infty) \\ \Delta_\infty &= L(\kappa_w, \kappa_n, \Delta_0, \Delta_\infty) = g^2 \int \mathcal{D}z \left[ \int \mathcal{D}x \phi(\kappa_w + \sqrt{\Delta_0 - \Delta_\infty} x + \sqrt{\Delta_\infty} z) \right]^2 + \kappa_n^2 + 1 \end{aligned} \tag{67}$$

833 where  $z$  here is a gaussian integration variable. We can solve these equations by simulating the  
834 following Langevin dynamical system.

$$\begin{aligned} x(t) &= \frac{\Delta_0(t)^2 - \Delta_\infty(t)^2}{2} \\ \Delta_0(t) &= \sqrt{2x(t) + \Delta_\infty(t)^2} \\ \dot{\kappa}_w(t) &= -\kappa_w(t) + F(\kappa_w(t), \kappa_n(t), \Delta_0(t), \Delta_\infty(t)) \\ \dot{\kappa}_n(t) &= -\kappa_n + G(\kappa_w(t), \kappa_n(t), \Delta_0(t), \Delta_\infty(t)) \\ \dot{x}(t) &= -x(t) + H(\kappa_w(t), \kappa_n(t), \Delta_0(t), \Delta_\infty(t)) \\ \dot{\Delta}_\infty(t) &= -\Delta_\infty(t) + L(\kappa_w(t), \kappa_n(t), \Delta_0(t), \Delta_\infty(t)) \end{aligned} \tag{68}$$

835 Then, the temporal variance, which is necessary for the Gaussian posterior conditioning example,  
836 is simply calculated via

$$\Delta_T = \Delta_0 - \Delta_\infty \quad (69)$$

837 **A.3 Supplementary Figures**

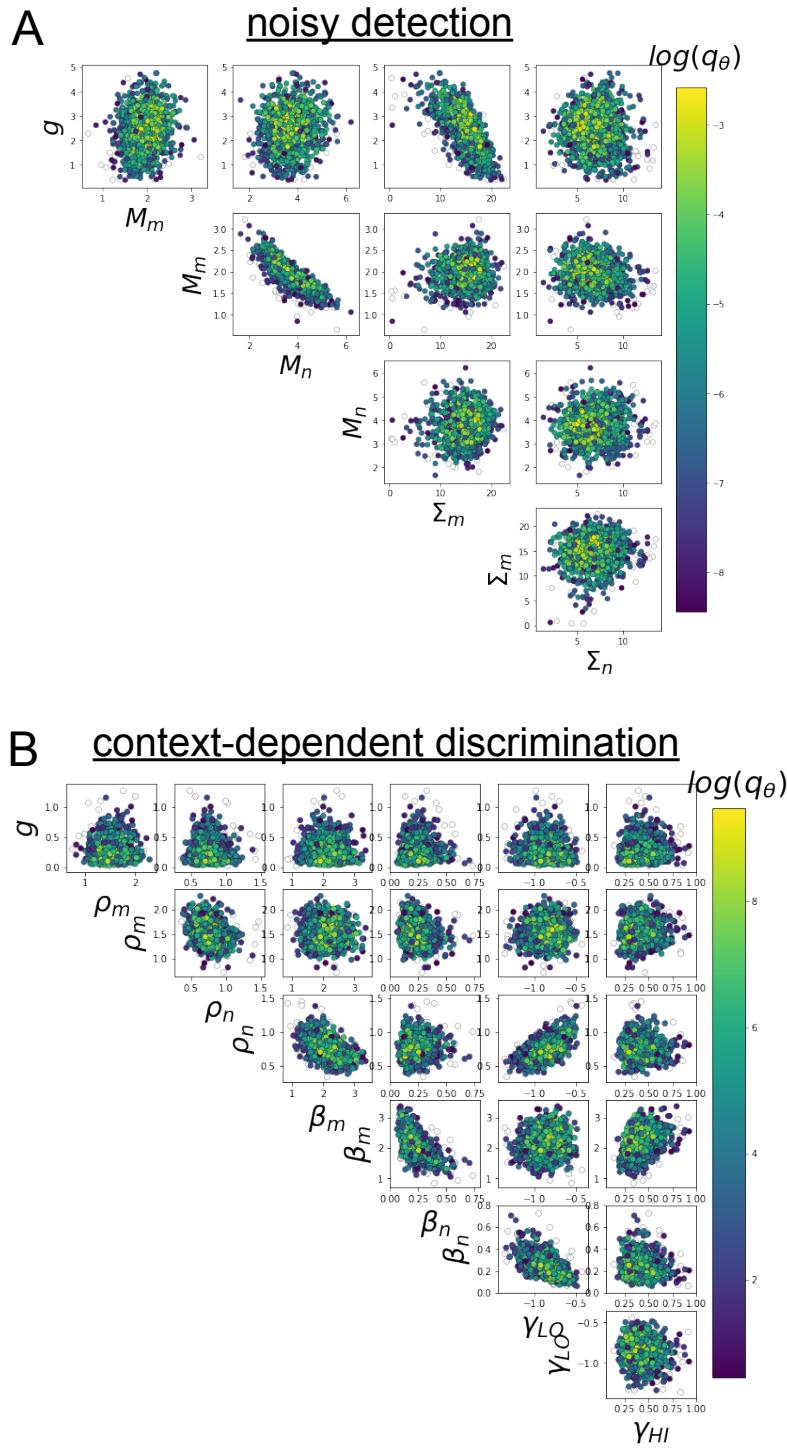


Fig. S1: A. EPI for rank-1 networks doing discrimination. B. EPI for rank-2 networks doing context-dependent discrimination.