

1 Training deep neural density estimators to 2 identify mechanistic models of neural dynamics

3 **Pedro J. Gonçalves^{1,2*}, Jan-Matthis Lueckmann^{1,2*}, Michael Deistler^{1*}, Marcel Nonnenmacher^{1,2},**
4 **Kaan Öcal^{2,3}, Giacomo Bassetto^{1,2}, Chaitanya Chintaluri⁴, William F. Podlaski⁴, Sara A. Haddad⁵,**
5 **Tim P. Vogels⁴, David S. Greenberg¹, Jakob H. Macke^{1,2}**

6 ¹Computational Neuroengineering, Department of Electrical and Computer Engineering, Technical
7 University of Munich, Germany; ²Max Planck Research Group Neural Systems Analysis, Center of
8 Advanced European Studies and Research (caesar), Bonn, Germany; ³Mathematical Institute, University of
9 Bonn, Bonn, Germany; ⁴Centre for Neural Circuits and Behaviour, University of Oxford; ⁵Max Planck
10 Institute for Brain Research, Frankfurt, Germany

12 Abstract

13 Mechanistic modeling in neuroscience aims to explain observed phenomena in terms of underlying causes. However,
14 determining which model parameters agree with complex and stochastic neural data presents a significant challenge.
15 We address this challenge with a machine learning tool which uses deep neural density estimators—trained using
16 model simulations—to carry out Bayesian inference and retrieve the full space of parameters compatible with raw
17 data or selected data features. Our method is scalable in parameters and data features, and can rapidly analyze new
18 data after initial training. We demonstrate the power and flexibility of our approach on receptive fields, ion channels,
19 and Hodgkin–Huxley models. We also characterize the space of circuit configurations giving rise to rhythmic activity in
20 the crustacean stomatogastric ganglion, and use these results to derive hypotheses for underlying compensation
21 mechanisms. Our approach will help close the gap between data-driven and theory-driven models of neural dynamics.

23 Introduction

24 New experimental technologies allow us to observe neurons, networks, brain regions and entire systems at un-
25 precedented scale and resolution, but using these data to understand how behavior arises from neural processes
26 remains a challenge. To test our understanding of a phenomenon, we often take to rebuilding it in the form of a
27 computational model that incorporates the mechanisms we believe to be at play, based on scientific knowledge,
28 intuition, and hypotheses about the components of a system and the laws governing their relationships. The goal of
29 such mechanistic models is to investigate whether a proposed mechanism can explain experimental data, uncover
30 details that may have been missed, inspire new experiments, and eventually provide insights into the inner workings
31 of an observed neural or behavioral phenomenon [1–4]. Examples for such a symbiotic relationship between model
32 and experiments range from the now classical work of Hodgkin and Huxley [5], to population models investigating
33 rules of connectivity, plasticity and network dynamics [6–10], network models of inter-area interactions [11, 12], and
34 models of decision making [13, 14].

35 A crucial step in building a model is adjusting its free parameters to be consistent with experimental observations.
36 This is essential both for investigating whether the model agrees with reality and for gaining insight into processes
37 which cannot be measured experimentally. For some models in neuroscience, it is possible to identify the relevant

*These authors contributed equally to this work

For correspondence: pedro.goncalves@caesar.de; jan-matthis.lueckmann@tum.de; michael.deistler@tum.de; macke@tum.de

38 parameter regimes from careful mathematical analysis of the model equations. But as the complexity of both neural
39 data and neural models increases, it becomes very difficult to find well-fitting parameters by inspection, and *automated*
40 identification of data-consistent parameters is required.

41 Furthermore, to understand how a model quantitatively explains data, it is necessary to find not only the *best*,
42 but *all* parameter settings consistent with experimental observations. This is especially important when modeling
43 neural data, where highly variable observations can lead to broad ranges of data-consistent parameters. Moreover,
44 many models in biology are inherently robust to some perturbations of parameters, but highly sensitive to others
45 [3, 15], e.g. because of processes such as homeostatic regulation. For these systems, identifying the full range of
46 data-consistent parameters can reveal how multiple distinct parameter settings give rise to the same model behavior
47 [7, 16, 17]. Yet, despite the clear benefits of mechanistic models in providing scientific insight, identifying their
48 parameters given data remains a challenging open problem that demands new algorithmic strategies.

49 The gold standard for automated parameter identification is *statistical inference*, which uses the likelihood $p(\mathbf{x}|\boldsymbol{\theta})$
50 to quantify the match between parameters $\boldsymbol{\theta}$ and data \mathbf{x} . Likelihoods can be derived for purely statistical models
51 commonly used in neuroscience [18–24], but are unavailable for most mechanistic models. Mechanistic models
52 are designed to reflect knowledge about biological mechanisms, and not necessarily to be amenable to efficient
53 inference: many mechanistic models are defined implicitly through stochastic computer simulations (e.g. a simulation
54 of a network of spiking neurons), and likelihood calculation would require the ability to integrate over all potential
55 paths through the simulator code. Similarly, a common goal of mechanistic modeling is to capture selected summary
56 features of the data (e.g. a certain firing rate, bursting behavior, etc...), *not* the full dataset in all its details. The same
57 feature (such as a particular average firing rate) can be produced by infinitely many realizations of the simulated
58 process (such as a time-series of membrane potential). This makes it impractical to compute likelihoods, as one would
59 have to average over all possible realizations which produce the same output.

60 Since the toolkit of statistical inference is inaccessible for mechanistic models, parameters are typically tuned
61 ad-hoc (often through laborious, and subjective, trial-and-error), or by computationally expensive parameter search: a
62 large set of models is generated, and grid search [25–27] or a genetic algorithm [28–31] is used to filter out simulations
63 which do not match the data. However, these approaches require the user to define a heuristic rejection criterion on
64 which simulations to keep (which can be challenging when observations have many dimensions or multiple units of
65 measurement), and typically end up discarding most simulations. Furthermore, they lack the advantages of statistical
66 inference, which provides principled approaches for handling variability, quantifying uncertainty, incorporating prior
67 knowledge and integrating multiple data sources. Approximate Bayesian Computation (ABC) [32–34] is a parameter-
68 search technique which aims to perform statistical inference, but still requires definition of a rejection criterion and
69 struggles in high-dimensional problems. Thus, computational neuroscientists face a dilemma: either create carefully
70 designed, highly interpretable mechanistic models (but rely on ad-hoc parameter tuning), or resort to purely statistical
71 models offering sophisticated parameter inference but limited mechanistic insight.

72 Here we propose a new approach using machine learning to combine the advantages of mechanistic and statistical
73 modeling. We present SNPE (Sequential Neural Posterior Estimation), a tool that rapidly identifies all mechanistic
74 model parameters consistent with observed experimental data (or summary features). SNPE builds on recent advances
75 in simulation-based Bayesian inference [35–38]: given observed experimental data (or summary features) \mathbf{x}_o , and a
76 mechanistic model with parameters $\boldsymbol{\theta}$, it expresses both prior knowledge and the range of data-compatible parameters
77 through probability distributions. SNPE returns a posterior distribution $p(\boldsymbol{\theta}|\mathbf{x}_o)$ which is high for parameters $\boldsymbol{\theta}$
78 consistent with both the data \mathbf{x}_o and prior knowledge, but approaches zero for $\boldsymbol{\theta}$ inconsistent with either (Fig. 1).

79 Similar to parameter search methods, SNPE uses simulations instead of likelihood calculations, but instead of
80 filtering out simulations, it uses *all* simulations to train a multi-layer artificial neural network to identify admissible
81 parameters (Fig. 1). By incorporating modern deep neural networks for conditional density estimation [39, 40], it can
82 capture the full *distribution* of parameters consistent with the data, even when this distribution has multiple peaks or
83 lies on curved manifolds. Critically, SNPE decouples the design of the model and design of the inference approach,
84 giving the investigator maximal flexibility to design and modify mechanistic models. Our method makes minimal
85 assumptions about the model or its implementation, and can e.g. also be applied to non-differentiable models, such
86 as networks of spiking neurons. Its only requirement is that one can run model simulations for different parameters,
87 and collect the resulting synthetic data or summary features of interest.

88 While the theoretical foundations of SNPE were developed and tested using simple inference problems on small

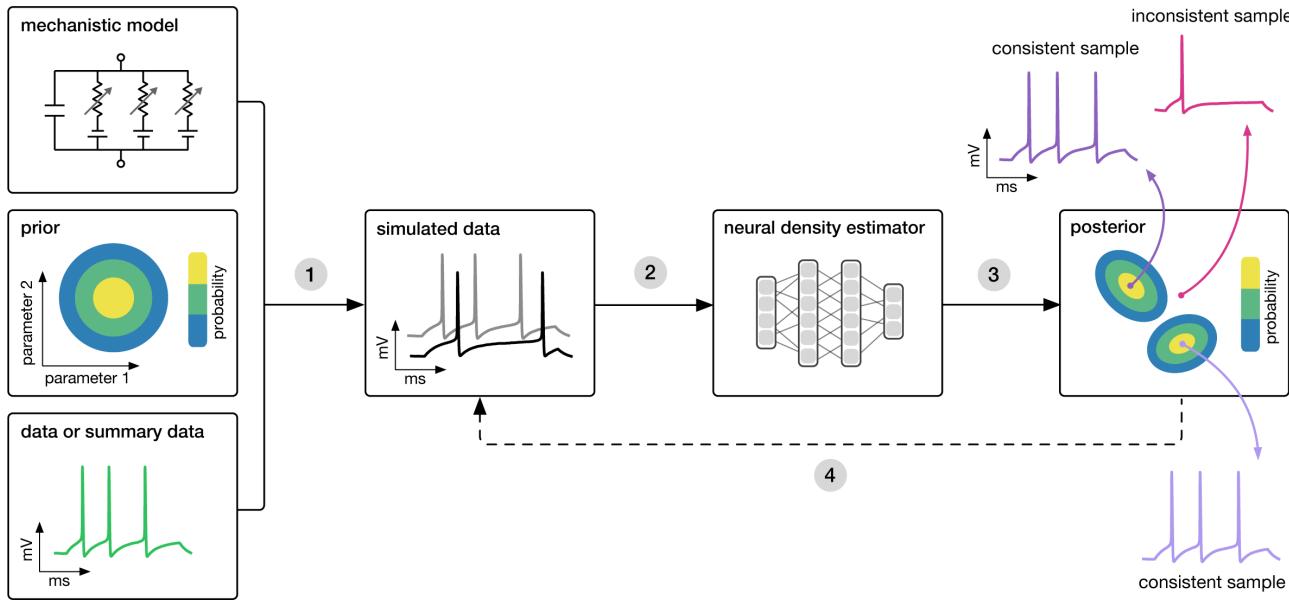


Figure 1. Goal: algorithmically identify mechanistic models which are consistent with data. Our algorithm (SNPE) takes three inputs: a candidate mechanistic model, prior knowledge or constraints on model parameters, and data (or summary statistics). SNPE proceeds by 1) sampling parameters from the prior and simulating synthetic datasets from these parameters, and 2) using a deep density estimation neural network to learn the (probabilistic) association between data (or data features) and underlying parameters, i.e. to learn statistical inference from simulated data. 3) This density estimation network is then applied to empirical data to derive the full space of parameters consistent with the data and the prior, i.e. the posterior distribution. High posterior probability is assigned to parameters which are consistent with both the data and the prior, low probability to inconsistent parameters. 4) If needed, an initial estimate of the posterior can be used to adaptively guide further simulations to produce data-consistent results.

models [35–37], here we show that SNPE can scale to complex mechanistic models in neuroscience, provide an accessible and powerful implementation, and develop validation and visualization techniques for exploring the derived posteriors. We illustrate SNPE using mechanistic models expressing key neuroscientific concepts: beginning with a simple neural encoding problem with a known solution, we progress to more complex data types, large datasets and many-parameter models inaccessible to previous methods. We estimate visual receptive fields using many data features, demonstrate rapid inference of ion channel properties from high-throughput voltage-clamp protocols, and show how Hodgkin–Huxley models are more tightly constrained by increasing numbers of data features. Finally, we showcase the power of SNPE by using it to identify the parameters of a network model which can explain an experimentally observed pyloric rhythm in the stomatogastric ganglion [7]—in contrast to previous approaches, SNPE allows us to search over the full space of both single-neuron and synaptic parameters, allowing us to study the geometry of the parameter space, as well as to provide new hypotheses for which compensation mechanisms might be at play.

Results

Estimating stimulus-selectivity in linear-nonlinear encoding models

We first illustrate SNPE on linear-nonlinear (LN) encoding models, a special case of generalized linear models (GLMs). These are simple, commonly used phenomenological models for which likelihood-based parameter estimation is feasible [41–46], and which can be used to validate the accuracy of our approach, before applying SNPE to more complex models for which the likelihood is unavailable. We will show that SNPE returns the correct posterior distribution over parameters, that it can cope with high-dimensional observation data, that it can recover multiple solutions to parameter inference problems, and that it is substantially more simulation efficient than conventional rejection-based ABC methods.

An LN model describes how a neuron’s firing rate is modulated by a sensory stimulus through a linear filter θ , often referred to as the *receptive field* [47, 48]. We first considered a model of a retinal ganglion cell (RGC) driven by full-field flicker (Fig. 2a). A statistic that is often used to characterize such a neuron is the *spike-triggered average* (STA) (Fig. 2a,

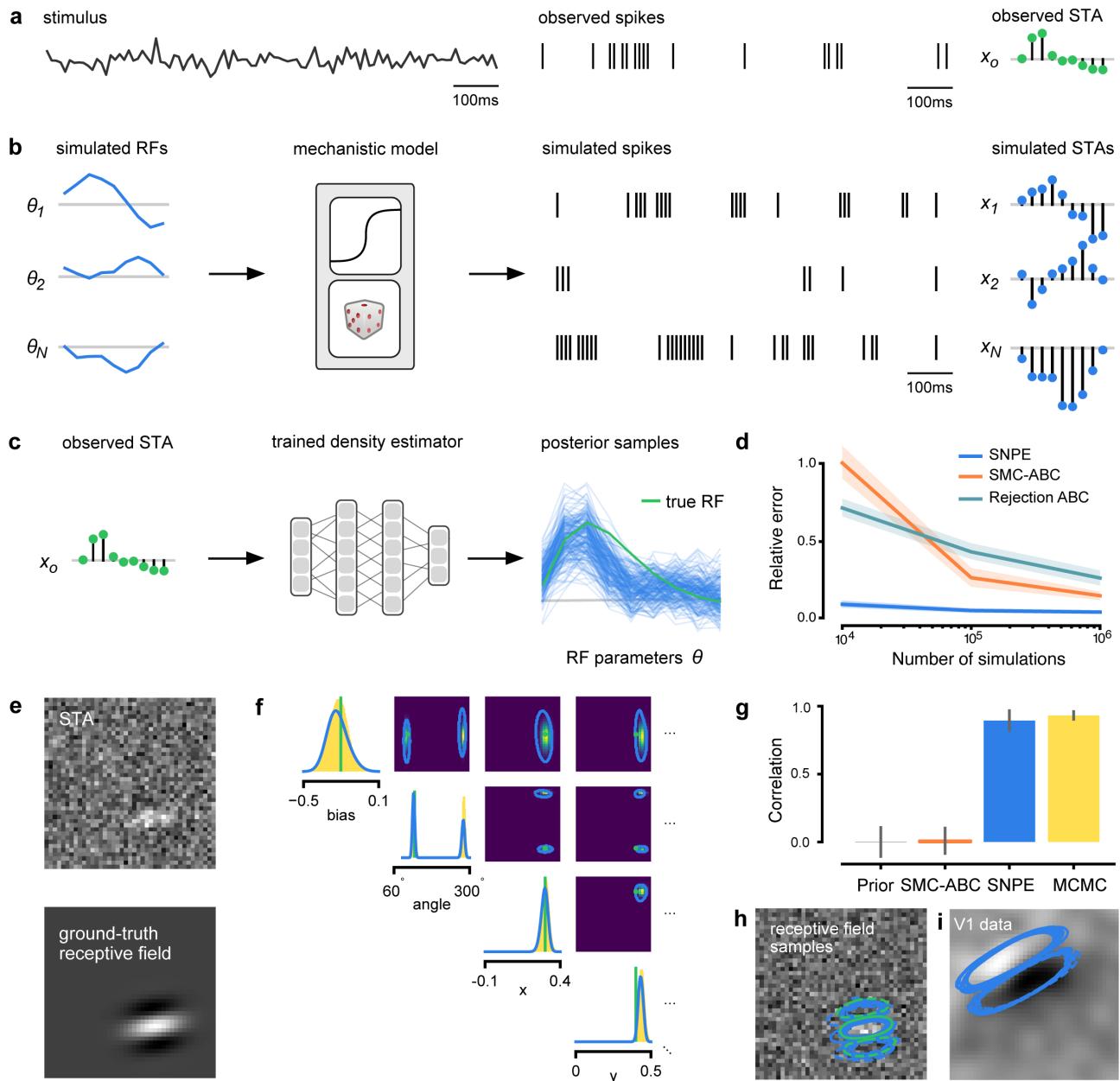


Figure 2. Estimating receptive fields in linear-nonlinear models of single neurons with statistical inference (a) Schematic of a time-varying stimulus, associated observed spike train and resulting spike-triggered average (STA). (b) SNPE proceeds by first randomly generating simulated receptive fields θ , and using the mechanistic model (here an LN model) to generate simulated spike trains and simulated STAs. (c) These simulated STAs and receptive fields are then used to train a deep neural density estimator to identify the distribution of receptive fields consistent with a given observed STA x_o . (d) Relative error in posterior estimation between SNPE and alternative methods (mean and 95%CI; 0 corresponds to perfect estimation, 1 to prior-level, details in Methods). (e) Example of spatial receptive field. We simulated responses and an STA of a LN-model with oriented receptive field. (f) We used SNPE to recover the distribution of receptive-field parameters. Univariate and pairwise marginals for four parameters of the spatial filter (MCMC, yellow histograms; SNPE, blue lines; full posterior in Supplementary Fig. 4). Non-identifiabilities of the Gabor parameterization lead to multimodal posteriors. (g) Average correlation (\pm SD) between ground-truth receptive field and receptive field samples from posteriors inferred with SMC-ABC, SNPE, and MCMC (which provides an upper bound given the inherent stochasticity of the data). (h) Posterior samples from SNPE posterior (SNPE, blue) compared to ground-truth receptive field (green; see panel (e)), overlaid on STA. (i) Posterior samples for V1 data; full posterior in Supplementary Fig. 5.

113 right). We therefore used the STA, as well as the firing rate of the neuron, as input x_o to SNPE. (Note that, in the limit of

114 infinite data, and for white noise stimuli, the STA will converge to the receptive field [42]—for finite, and non-white data,
115 the two will in general be different.) Starting with random receptive fields θ , we generated synthetic spike trains and
116 calculated STAs from them (Fig. 2b). We then trained a neural conditional density estimator to recover the receptive
117 fields from the STAs and firing rates (Fig. 2c). This allowed us to estimate the posterior distribution over receptive fields,
118 i.e. to estimate which receptive fields are consistent with the data (and prior) (Fig. 2c). For LN models, likelihood-based
119 inference is possible, allowing us to validate the SNPE posterior by comparing it to a reference posterior obtained
120 via Markov Chain Monte Carlo (MCMC) sampling [45, 46]. We found that SNPE accurately estimates the posterior
121 distribution (Supplementary Fig. 1 and Supplementary Fig. 2), and substantially outperforms Sequential Monte Carlo
122 (SMC) ABC methods [34, 49] (Fig. 2d).

123 As a more challenging problem, we inferred the receptive field of a neuron in primary visual cortex (V1) [50, 51].
124 Using a model composed of a bias (related to the spontaneous firing rate) and a Gabor function with 8 parameters
125 [52] describing the receptive field's location, shape and strength, we simulated responses to 5-minute random noise
126 movies of 41×41 pixels, such that the STA is high-dimensional, with a total of 1681 dimensions (Fig. 2e). This problem
127 admits multiple solutions (as e.g. rotating the receptive field by 180°). As a result, the posterior distribution has
128 multiple peaks ('modes'). Starting from a simulation result x_o with known parameters, we used SNPE to estimate the
129 posterior distribution $p(\theta|x_o)$. To deal with the high-dimensional data x_o in this problem, we used a convolutional
130 neural network (CNN), as this architecture excels at learning relevant features from image data [53, 54]. To deal with
131 the multiple peaks in the posterior, we fed the CNN's output into a mixture density network (MDN) [55], which can
132 learn to assign probability distributions with multiple peaks as a function of its inputs (details in Methods). Using this
133 strategy, SNPE was able to infer a posterior distribution that tightly enclosed the ground truth simulation parameters
134 which generated the original simulated data x_o , and matched a reference MCMC posterior (Fig. 2f, posterior over
135 all parameters in Supplementary Fig. 4). For this challenging estimation problem with high-dimensional summary
136 features, an SMC ABC algorithm with the same simulation-budget failed to identify the correct receptive fields (Fig. 2g)
137 and posterior distributions (Supplementary Fig. 3). We also applied this approach to electrophysiological data from a
138 V1 cell [51], identifying a sine-shaped Gabor receptive field consistent with the original spike-triggered average (Fig. 2i);
139 posterior distribution in Supplementary Fig. 5).

140 **Functional diversity of ion channels: efficient high-throughput inference**

141 We next show how SNPE can be efficiently applied to estimation problems in which we want to identify a large number
142 of models for different observations in a database. We considered a flexible model of ion channels [57], which we
143 here refer to as the *Omnimodel*. This model uses 8 parameters to describe how the dynamics of currents through
144 non-inactivating potassium channels depend on membrane voltage (Fig. 3a). For various choices of its parameters θ ,
145 it can capture 350 specific models in publications describing this channel type, cataloged in the IonChannelGenealogy
146 (ICG) database [56]. We aimed to identify these ion channel parameters θ for each ICG model, based on 11 features
147 of the model's response to a sequence of 5 voltage clamp protocols, resulting in a total of 55 different characteristic
148 features per model (Fig. 3b, see Methods for details).

149 Because this model's output is a typical format for functional characterization of ion channels both in simulations
150 [56] and in high-throughput electrophysiological experiments [58–60], the ability to rapidly infer different parameters
151 for many separate experiments is advantageous. Existing approaches for fitting deterministic models based on
152 numerical optimization [57, 60] must repeat all computations anew for a new experiment or data point (Fig. 3c).
153 However, for SNPE the only heavy computational tasks are carrying out simulations to generate training data, and
154 training the neural network. We therefore reasoned that by training a network once using a large number of
155 simulations, we could subsequently carry out rapid 'amortized' parameter inference on new data using a single pass
156 through the network (Fig. 3d) [61, 62]. To test this idea, we used SNPE to train a neural network to infer the posterior
157 from any data x . To generate training data, we carried out 1 million Omnimodel simulations, with parameters randomly
158 chosen across ranges large enough to capture the models in the ICG database [56]. SNPE was run using a single round,
159 i.e. it learned to perform inference for all data from the prior (rather than a specific observed datum). Generating these
160 simulations took around 1000 CPU-hours and training the network 150 CPU-hours, but afterwards a full posterior
161 distribution could be inferred for new data in less than 10 ms.

162 As a first test, SNPE was run on simulation data, generated by a previously published model of a non-inactivating
163 potassium channel [63] (Fig. 3b). Simulations of the Omnimodel using parameter sets sampled from the obtained

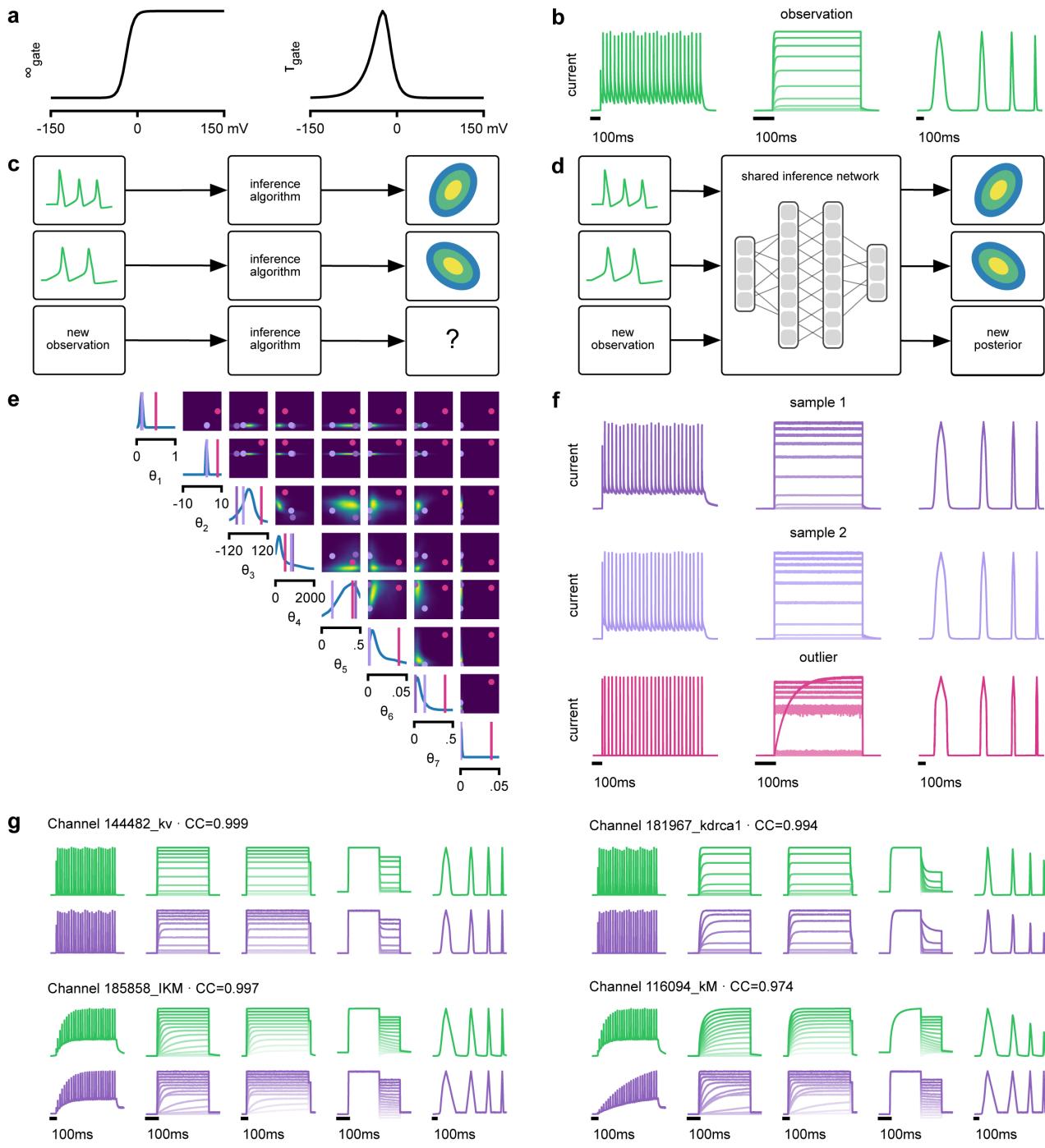


Figure 3. Inference on a database of ion-channel models. (a) We perform inference over the parameters of non-inactivating potassium channel models. Channel kinetics are described by steady-state activation curves, ∞_{gate} , and time-constant curves, τ_{gate} . (b) Observation generated from a channel model from ICG database: normalized current responses to three (out of five) voltage-clamp protocols (action potentials, activation, and ramping). Details in [56]. (c) Classical approach to parameter identification: inference is optimized on each datum separately, requiring new computations for each new datum. (d) Amortized inference: an inference network is learned which can be applied to multiple data, enabling rapid inference on new data. (e) Posterior distribution over eight model parameters, θ_1 to θ_8 . (f) Traces obtained by sampling from the posterior in (e). Purple: traces sampled from posterior, i.e. with high posterior probability. Magenta: trace from parameters with low probability. (g) Observations (green) and traces generated by posterior samples (purple) for four models from the database.

164 posterior distribution (Fig. 3e) closely resembled the input data on which the SNPE-based inference had been carried
165 out, while simulations using ‘outlier’ parameter sets with low probability under the posterior generated current
166 responses that were markedly different from the data x_o (Fig. 3f). Taking advantage of SNPE’s capability for rapid
167 amortized inference, we further evaluated its performance on all 350 non-inactivating potassium channel models
168 in ICG. In each case, we carried out a simulation to generate initial data from the original ICG model, used SNPE to
169 calculate the posterior given the Omnimodel, and then generated a new simulation x using parameters sampled from
170 the posterior (Fig. 3f). This resulted in high correlation between the original ICG model response and the Omnimodel
171 response, in every case (>0.98 for more than 90% of models, see Supplementary Fig. 6). However, this approach was
172 not able to capture all traces perfectly, as e.g. it failed to capture the shape of the onset of the bottom right model in
173 Fig. 3g. Additional analysis of this example revealed that this example is not a failure of SNPE, but rather a limitation
174 of the Omnimodel. Thus, SNPE can be used to reveal limitations of candidate models and aid the development of
175 more verisimilar mechanistic models.

176 Calculating the posterior for all 350 ICG models took only a few seconds, and was fully automated, i.e. did not
177 require user interactions. These results show how SNPE allows fast and accurate identification of biophysical model
178 parameters on new data, and how SNPE can be deployed for applications requiring rapid automated inference,
179 such as high-throughput screening-assays, closed-loop paradigms (e.g. for adaptive experimental manipulations or
180 stimulus-selection), or interactive software tools.

181 **Hodgkin-Huxley model: stronger constraints from additional data features**

182 The Hodgkin-Huxley (HH) model [5] of action potential generation through ion channel dynamics is a highly influential
183 mechanistic model in neuroscience. A number of algorithms have been proposed for fitting HH models to electrophys-
184 iological data [25, 30, 31, 64–67], but [with the exception of 68] these approaches do not attempt to estimate the full
185 posterior. Given the central importance of the HH model in neuroscience, we sought to test how SNPE would cope
186 with this challenging non-linear model.

187 As previous approaches for HH models concentrated on reproducing specified features [e.g. the number of spikes,
188 65], we also sought to determine how various features provide different constraints. We considered the problem of
189 inferring 8 biophysical parameters in a HH single-compartment model, describing voltage-dependent sodium and
190 potassium conductances and other intrinsic membrane properties (Fig. 4a, left). We simulated the neuron’s voltage
191 response to the injection of a square wave of depolarizing current, and defined the model output x used for inference
192 as the number of evoked action potentials along with 6 additional features of the voltage response (Fig. 4a, right,
193 details in Methods). We first applied SNPE to observed data x_o created by simulation from the model, calculating
194 the posterior distribution using all 7 features in the observed data (Fig. 4b). The posterior contained the ground
195 truth parameters in a high probability-region, as in previous applications, indicating the consistency of parameter
196 identification. The variance of the posterior was narrower for some parameters than for others, indicating that the
197 7 data features constrain some parameters strongly (such as the potassium conductance), but others only weakly
198 (such as the adaptation time constant). Additional simulations with parameters sampled from the posterior closely
199 resembled the observed data x_o , in terms of both the raw membrane voltage over time and the 7 data features (Fig. 4c,
200 purple and green). Parameters with low posterior probability (outliers) generated simulations that markedly differed
201 from x_o (Fig. 4c, magenta).

202 Genetic algorithms are commonly used to fit parameters of deterministic biophysical models [28, 29, 31, 69].
203 While genetic algorithms can also return multiple data-compatible parameters, they do not perform inference (i.e.
204 find the posterior distribution), and their outputs depend strongly on user-defined goodness-of-fit criteria. When
205 comparing a state-of-the-art genetic algorithm [Indicator Based Evolutionary Algorithm, IBEA, 31, 70, 71] to SNPE,
206 we found that the parameter-settings favoured by IBEA produced simulations whose summary features were as
207 similar to the observed data as those obtained by SNPE high-probability samples (Supplementary Fig. 9). However,
208 high-scoring IBEA parameters were concentrated in small regions of the posterior, i.e. IBEA did not identify the full
209 space of data-compatible models.

210 To investigate how individual data features constrain parameters, we compared SNPE-estimated posteriors based
211 1) solely on the spike count, 2) on the spike count and 3 voltage-features, or 3) on all 7 features of x_o . As more features
212 were taken into account, the posterior became narrower and centered more closely on the ground truth parameters
213 (Fig. 4d, Supplementary Fig. 7). Posterior simulations matched the observed data only in those features that had been

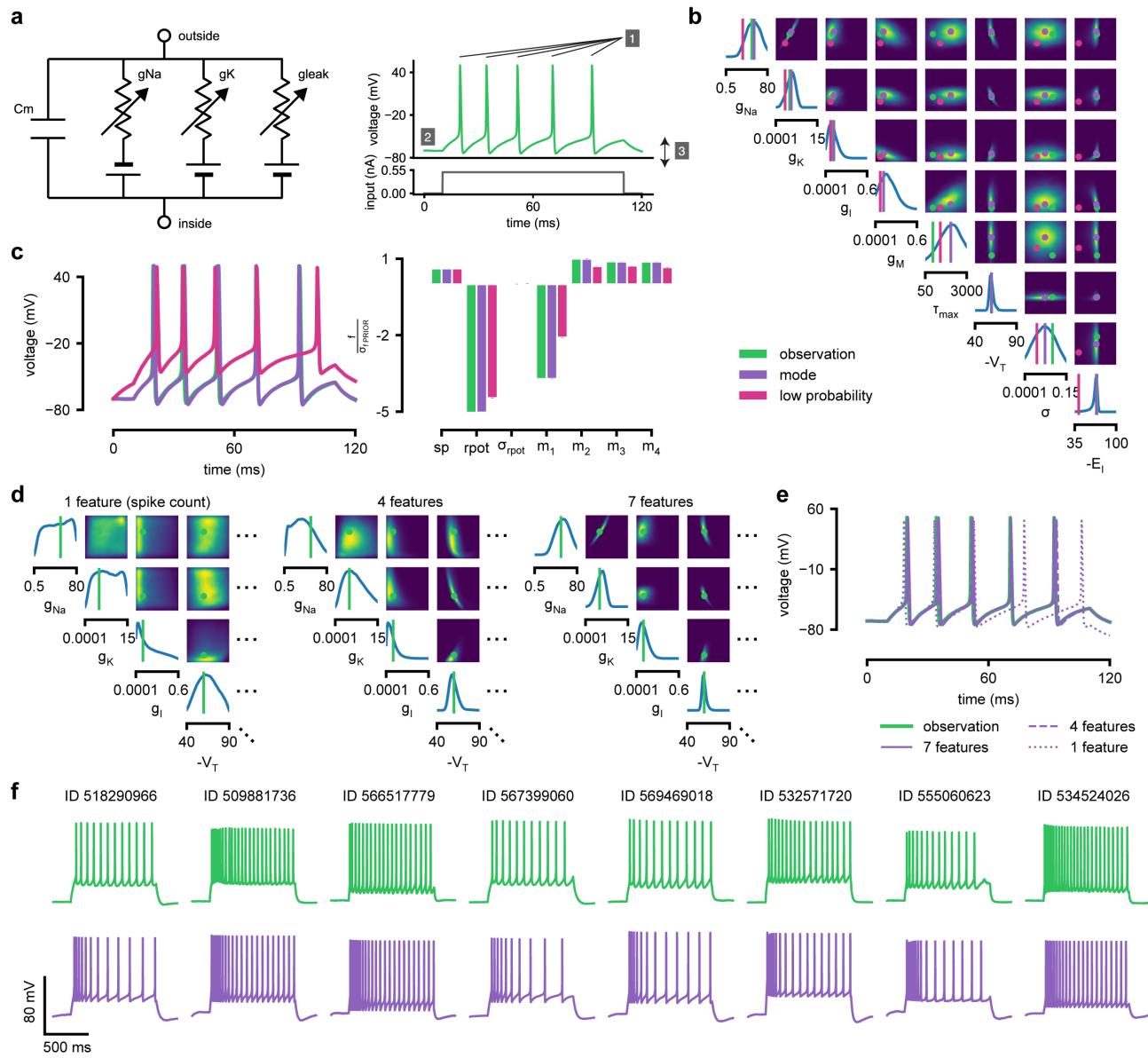


Figure 4. Inference for single compartment Hodgkin-Huxley model. (a) Circuit diagram describing the Hodgkin-Huxley model (left), and simulated voltage-trace given a current input (right). 3 out of 7 voltage features are depicted: (1) number of spikes, (2) mean resting potential and (3) standard deviation of the pre-stimulus resting potential. (b) Inferred posterior for 8 parameters given 7 voltage features. (c) Traces (left) and associated features f (right) for the desired output (observation), the mode of the inferred posterior, and a sample with low posterior probability. The voltage features are: number of spikes sp , mean resting potential $rpot$, standard deviation of the resting potential σ_{rpot} , and the first 4 voltage moments, mean m_1 , standard deviation m_2 , skewness m_3 and kurtosis m_4 . Each feature is normalized by $\sigma_f \text{ PRIOR}$, the standard deviation of the respective feature of simulations sampled from the prior. (d) Partial view of the inferred posteriors (4 out of 8 parameters) given 1, 4 and 7 features (full posteriors over 8 parameters in Supplementary Fig. 7). (e) Traces for posterior modes given 1, 4 and 7 features. Increasing the number of features leads to posterior traces that are closer to the observed data. (f) Observations from Allen Cell Types Database (green) and corresponding mode samples (purple). Posteriors in Supplementary Fig. 8.

used for inference (e.g. applying SNPE to spike counts alone identified parameters that generated the correct number of spikes, but for which spike timing and subthreshold voltage time course were off, Fig. 4e). For some parameters, such as the potassium conductance, providing more data features brought the peak of the posterior (the *posterior mode*) closer to the ground truth and also decreased uncertainty. For other parameters, such as V_T , a parameter adjusting the spike threshold [65], the peak of the posterior was already close to the correct value with spike counts

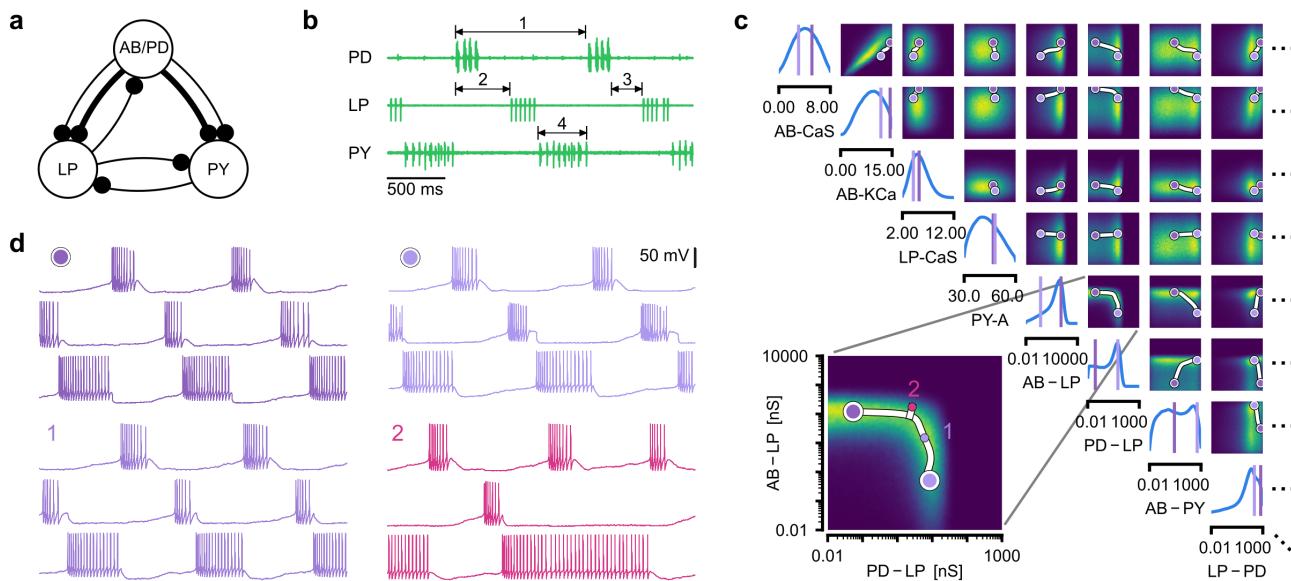


Figure 5. Identifying network models underlying an experimentally observed pyloric rhythm in the crustacean stomatogastric ganglion. (a) Simplified circuit diagram of the pyloric network from the stomatogastric ganglion. Thin connections are fast glutamatergic, thick connections are slow cholinergic. (b) Extracellular recordings from nerves of pyloric motor neurons of the crab *Cancer borealis* [74]. Numbers indicate some of the used summary features, namely cycle period (1), phase delays (2), phase gaps (3), and burst durations (4) (see Methods for details). (c) Posterior over 24 membrane and 7 synaptic conductances given the experimental observation shown in panel b (8 parameters shown, full posterior in Supplementary Fig. 10). Inset: magnified marginal posterior for the synaptic strengths AB to LP neuron vs. PD to LP neuron. (d) Identifying directions of slippiness and stiffness. Two samples from the posterior both show similar network activity as the experimental observation (top left and top right), but have very different parameters (purple dots in panel c). Along the high-probability path between these samples, network activity is preserved (trace 1). When perturbing the parameters orthogonally off the path, network activity changes abruptly and becomes non-pyloric (trace 2).

219 alone, but adding additional features reduced uncertainty. While SNPE can be used to study the effect of additional
220 data features in reducing parameter uncertainty, this would not be the case for methods that only return a single
221 best-guess estimate of parameters. These results show that SNPE can reveal how information from multiple data
222 features imposes collective constraints on channel and membrane properties in the HH model.

223 We also inferred HH parameters for 8 *in vitro* recordings from the Allen Cell Types database using the same current-
224 clamp stimulation protocol as in our model [72, 73] (Fig. 4f, Supplementary Fig. 8). In each case, simulations based
225 on the SNPE-inferred posterior closely resembled the original data (Fig. 4f). We note that while inferred parameters
226 differed across recordings, some parameters (the spike threshold, the density of sodium channels, the membrane
227 reversal potential and the density of potassium channels) were consistently more strongly constrained than others
228 (the intrinsic neural noise, the adaptation time constant, the density of slow voltage-dependent channels and the leak
229 conductance) (Supplementary Fig. 8). Overall, these results suggest that the electrophysiological responses measured
230 by this current-clamp protocol can be approximated by a single-compartment HH model, and that SNPE can identify
231 the admissible parameters.

232 Crustacean stomatogastric ganglion: sensitivity to perturbations

233 We next aimed to demonstrate how the full posterior distribution obtained with SNPE can lead to novel scientific
234 insights. To do so, we used the pyloric network of the stomatogastric ganglion (STG) of the crab *Cancer borealis*, a
235 well-characterized neural circuit producing rhythmic activity. In this circuit, similar network activity can arise from
236 vastly different sets of membrane and synaptic conductances [7]. We first investigated whether data-consistent sets
237 of membrane and synaptic conductances are connected in parameter space, as has been demonstrated for single
238 neurons [75], and, second, which compensation mechanisms between parameters of this circuit allow the neural
239 system to maintain its activity despite parameter variations. While this model has been studied extensively, answering
240 these questions requires characterizing higher-dimensional parameter spaces than those accessed previously. We

241 demonstrate how SNPE can be used to identify the posterior distribution over both membrane and synaptic con-
242 ductances of the STG (31 parameters total) and how the full posterior distribution can be used to study the above
243 questions at the circuit level.

244 For some biological systems, multiple parameter sets give rise to the same system behavior [7, 17, 76–79]. In
245 particular, neural systems can be robust to specific perturbations of parameters [79–81], yet highly sensitive to others,
246 properties referred to as *sloppiness* and *stiffness* [3, 15, 82, 83]. We studied how perturbations affect model output
247 using a model [7] and data [74] of the pyloric rhythm in the crustacean stomatogastric ganglion (STG). This model
248 describes a triphasic motor pattern generated by a well-characterized circuit (Fig. 5a). The circuit consists of two
249 electrically coupled pacemaker neurons (anterior burster and pyloric dilator, AB/PD), modeled as a single neuron, as
250 well as two types of follower neurons (lateral pyloric (LP) and pyloric (PY)), all connected through inhibitory synapses
251 (details in Methods). Eight membrane conductances are included for each modeled neuron, along with 7 synaptic
252 conductances, for a total of 31 parameters. This model has been used to demonstrate that virtually indistinguishable
253 activity can arise from vastly different membrane and synaptic conductances in the STG [7, 17].

254 We applied SNPE to an extracellular recording from the STG of the crab *Cancer borealis* [74] which exhibited pyloric
255 activity (Fig. 5b), and inferred the posterior distribution over all 31 parameters based on 18 salient features of the
256 voltage traces, including cycle period, phase delays, phase gaps, and burst durations (features in Fig. 5B, posterior
257 in Fig. 5c, posterior over all parameters in Supplementary Fig. 10, details in Methods). Consistent with previous
258 reports, the posterior distribution has high probability over extended value ranges for many membrane and synaptic
259 conductances. To verify that parameter settings across these extended ranges are indeed capable of generating the
260 experimentally observed network activity, we sampled two sets of membrane and synaptic conductances from the
261 posterior distribution. These two samples have widely disparate parameters from each other (Fig. 5c, purple dots,
262 details in Methods), but both exhibit activity highly similar to the experimental observation (Fig. 5d, top left and top
263 right).

264 We then investigated the geometry of the parameter space producing these rhythms [16, 17]. First, we wanted to
265 identify directions of sloppiness, and we were interested in whether parameter settings producing pyloric rhythms
266 form a single connected region, as has been shown for single neurons [75], or whether they lie on separate ‘islands.’
267 Starting from the two above parameter settings showing similar activity, we examined whether they were connected
268 by searching for a path through parameter space along which pyloric activity was maintained. To do this, we
269 algorithmically identified a path lying only in regions of high posterior probability (Fig. 5c, white, details in Methods).
270 Along the path, network output was tightly preserved, despite a substantial variation of the parameters (voltage trace
271 1 in Fig. 5d, Supplementary Fig. 11a,c). Second, we inspected directions of stiffness by perturbing parameters off
272 the path. We applied perturbations that yield maximal drops in posterior probability (see Methods for details), and
273 found that the network quickly produced non-pyloric activity (voltage trace 2, Fig. 5d) [82]. In identifying these paths
274 and perturbations, we exploited the fact that SNPE provides a differentiable estimate of the posterior, as opposed to
275 parameter search methods which provide only discrete samples.

276 Overall, these results show that the pyloric network can be robust to specific perturbations in parameter space,
277 but sensitive to others, and that one can interpolate between disparate solutions while preserving network activity.
278 This analysis demonstrates the flexibility of SNPE in capturing complex posterior distributions, and shows how the
279 differentiable posterior can be used to study directions of sloppiness and stiffness.

280 **Predicting compensation mechanisms from posterior distributions**

281 Experimental and computational studies have shown that stable neural activity can be maintained despite variable
282 circuit parameters [7, 87, 88]. This behavior can emerge from two sources [87]: either, the variation of a certain
283 parameter barely influences network activity at all, or alternatively, variations of several parameters influence network
284 activity, but their effects compensate for one another. Here, we investigated these possibilities by using the posterior
285 distribution over membrane and synaptic conductances of the STG.

286 We begin by drawing samples from the posterior and inspecting their pairwise histograms (i.e. the pairwise
287 marginals, Fig. 6a, posterior over all parameters in Supplementary Fig. 10). Consistent with previously reported results
288 [89], many parameters seem only weakly constrained and only weakly correlated (Fig. 6b). However, this observation
289 does not imply that the parameters of the network do not have to be finely tuned: pairwise marginals are averages
290 over many network configurations, where all other parameters may take on diverse values, which could disguise that

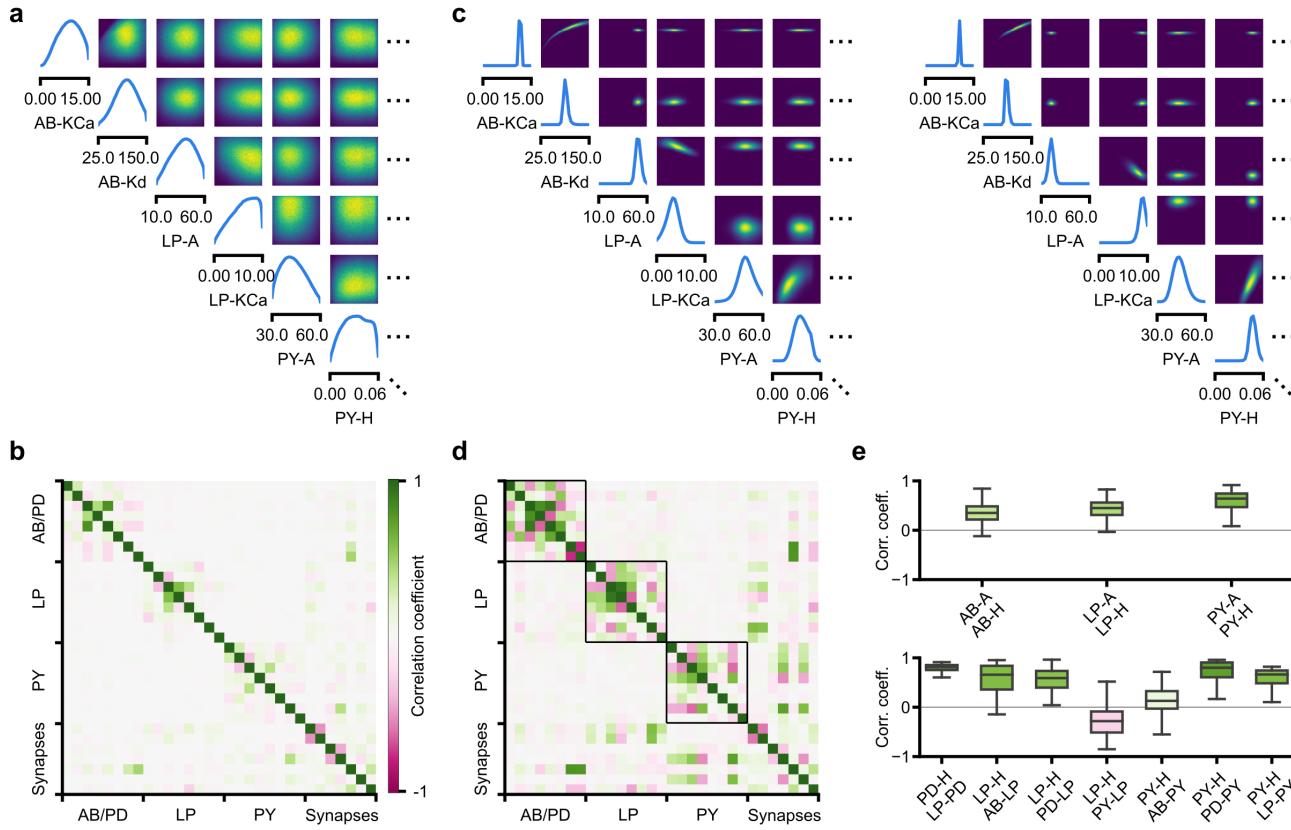


Figure 6. Predicting compensation mechanisms in the stomatogastric ganglion. (a) Inferred posterior. We show a subset of parameters which are weakly constrained (full posterior in Supplementary Fig. 10). Pyloric activity can emerge from a wide range of maximal membrane conductances, as the 1D and 2D posterior marginals cover almost the entire extent of the prior. (b) Correlation matrix, based on the samples shown in panel a. Almost all correlations are weak. Ordering of membrane and synaptic conductances as in Supplementary Fig. 10. (c) Conditional distributions given a particular circuit configuration: for the plots on the diagonal, we keep all but one parameter fixed. For plots above the diagonal, we keep all but two parameters fixed. The remaining parameter(s) are narrowly tuned, tuning across parameters is often highly correlated. When conditioning on a different parameter setting (right plot), the conditional posteriors change, but correlations are often maintained. (d) Conditional correlation matrix, averaged over 500 conditional distributions like the ones shown in panel c. Black squares highlight parameter-pairs within the same model neuron. (e) Consistency with experimental observations. Top: maximal conductance of the fast transient potassium current and the maximal conductance of the hyperpolarization current are positively correlated for all three neurons. This has also been experimentally observed in the PD and the LP neuron [84]. Bottom: the maximal conductance of the hyperpolarization current of the postsynaptic neuron can compensate the strength of the synaptic input, as experimentally observed in the PD and the LP neuron [85, 86]. The boxplots indicate the maximum, 75% quantile, median, 25% quantile, and minimum across 500 conditional correlations for different parameter pairs. Face color indicates mean correlation using the colorbar shown in panel b.

each individual configuration is finely tuned. Indeed, when we sampled parameters independently from their posterior histograms, the resulting circuit configurations rarely produced pyloric activity, indicating that parameters have to be tuned relative to each other (Supplementary Fig. 12). This analysis also illustrates that the (common) approach of independently setting parameters can be problematic: although each parameter individually is in a realistic range, the network as a whole is not [90]. Finally, it shows the importance of identifying the full posterior distribution, which is far more informative than just finding individual parameters and assigning error bars.

In order to investigate the need for tuning between pairs of parameters, we held all but two parameters constant at a given consistent circuit configuration (sampled from the posterior), and observed the network activity across different values of the remaining pair of parameters. We can do so by calculating the conditional posterior distribution (details in Methods), and do not have to generate additional simulations (as would be required by parameter search methods). Doing so has a simple interpretation: when all but two parameters are fixed, what values of the remaining two parameters can then lead to the desired network activity? We found that the desired pattern of pyloric activity

303 can emerge only from narrowly tuned and often highly correlated combinations of the remaining two parameters,
304 showing how these parameters can compensate for one another (Fig. 6c). When repeating this analysis across multiple
305 network configurations, we found that these ‘conditional correlations’ are often preserved (Fig. 6c, left and right). This
306 demonstrates that pairs of parameters can compensate for each other in a similar way, independently of the values
307 taken by other parameters. This observation about compensation could be interpreted as an instance of modularity, a
308 widespread underlying principle of biological robustness [91].

309 We calculated conditional correlations for each parameter pair using 500 different circuit configurations sampled
310 from the posterior (Fig. 6d). Compared to correlations based on the pairwise marginals (Fig. 6b), these conditional
311 correlations were substantially stronger. They were particularly strong across membrane conductances of the same
312 neuron, but primarily weak across different neurons (black boxes in Fig. 6d).

313 Finally, we tested whether the conditional correlations were in line with experimental observations. For the PD
314 and the LP neuron, it has been reported that overexpression of the fast transient potassium current (I_A) leads to a
315 compensating increase of the hyperpolarization current (I_H), suggesting a positive correlation between these two
316 currents [84, 92]. These results are qualitatively consistent with the positive conditional correlations between the
317 maximal conductances of I_A and I_H for all three model neurons (Fig. 6e top). In addition, using the dynamic clamp, it
318 has been shown that diverse combinations of the synaptic input strength and the maximal conductance of I_H lead
319 to similar activity in the LP and the PD neuron [85, 86]. Consistent with these findings, the non-zero conditional
320 correlations reveal that there can indeed be compensation mechanisms between the synaptic strength and the
321 maximal conductance of I_H of the postsynaptic neuron (Fig. 6e bottom).

322 Overall, we showed how SNPE can be used to study parameter dependencies, and how the posterior distribution
323 can be used to efficiently explore potential compensation mechanisms. We found that our method can predict
324 compensation mechanisms which are qualitatively consistent with experimental studies. We emphasize that these
325 findings would not have been possible with a direct grid-search over all parameters: defining a grid in a 31-dimensional
326 parameter space would require more than $2^{31} > 2$ billion simulations, even if one were to use the coarsest-possible
327 grid with only 2 values per dimension.

328 Discussion

329 How can we build models which give insights into the causal mechanisms underlying neural or behavioral dynamics?
330 The cycle of building mechanistic models, generating predictions, comparing them to empirical data, and rejecting
331 or refining models has been of crucial importance in the empirical sciences. However, a key challenge has been the
332 difficulty of identifying mechanistic models which can quantitatively capture observed phenomena. We suggest that a
333 generally applicable tool to constrain mechanistic models by data would expedite progress in neuroscience. While
334 many considerations should go into designing a model that is appropriate for a given question and level of description
335 [2, 3, 93, 94], the question of whether and how one can perform statistical inference should not compromise model
336 design. In our tool, SNPE, the process of model building and parameter inference are entirely decoupled. SNPE can be
337 applied to *any* simulation-based model (requiring neither model nor summary features to be differentiable) and gives
338 full flexibility on defining a prior. We illustrated the power of our approach on a diverse set of applications, highlighting
339 the potential of SNPE to rapidly identify data-compatible mechanistic models, to investigate which data-features
340 effectively constrain parameters, and to reveal shortcomings of candidate-models.

341 Finally, we used a model of the stomatogastric ganglion to show how SNPE can identify complex, high-dimensional
342 parameter landscapes of neural systems. We analyzed the geometrical structure of the parameter landscape and
343 confirmed that circuit configurations need to be finely tuned, even if individual parameters can take on a broad range
344 of values. We showed that different configurations are connected in parameter space, and provided hypotheses for
345 compensation mechanisms. These analyses were made possible by SNPE’s ability to estimate full parameter posteriors,
346 rather than just constraints on individual parameters, as is common in many statistical parameter-identification
347 approaches.

348 Related work

349 SNPE builds on recent advances in machine learning, and in particular in density-estimation approaches to likelihood-
350 free inference [35–37, 95, 96], reviewed in [38]. We here scaled these approaches to canonical mechanistic models
351 of neural dynamics, and provided methods and software-tools for inference, visualization, and analysis of the

352 resulting posteriors (e.g. the high-probability paths and conditional correlations presented here). The idea of learning
353 inference networks on simulated data can be traced back to *regression-adjustment* methods in ABC [32, 97]. [35]
354 first proposed to use expressive conditional density estimators in the form of deep neural networks [40, 55], and
355 to optimize them sequentially over multiple rounds with cost-functions derived from Bayesian inference principles.
356 Compared to commonly used rejection-based ABC methods [98, 99], such as MCMC-ABC [33], SMC-ABC [34, 100],
357 Bayesian-Optimization ABC [101], or ensemble methods [102, 103], SNPE approaches do not require one to define
358 a distance function in data space. In addition, by leveraging the ability of neural networks to learn informative
359 features, they enable scaling to problems with high-dimensional observations, as are common in neuroscience
360 and other fields in biology. We have illustrated this capability in the context of receptive field estimation, where a
361 convolutional neural network extracts summary features from a 1681 dimensional spike-triggered average. Alternative
362 likelihood-free approaches include *synthetic likelihood* methods [104–110], moment-based approximations of the
363 posterior [111, 112], inference compilation [113, 114], and density-ratio estimation [115]. For some mechanistic
364 models in neuroscience (e.g. for integrate-and-fire neurons), likelihoods can be computed via stochastic numerical
365 approximations [66, 116, 117] or model-specific analytical approaches [64, 118–121].

366 Our approach is already finding its first applications in neuroscience—for example, [122] have used a variant of
367 SNPE to constrain biophysical models of retinal neurons, with the goal of optimizing stimulation approaches for
368 neuroprosthetics. Concurrently with our work, [123] developed an alternative approach to parameter identification
369 for mechanistic models, and showed how it can be used to characterize neural population models which exhibit
370 specific emergent computational properties. Both studies differ in their methodology and domain of applicability
371 (see descriptions of underlying algorithms in our [36, 37] and their [124] prior work), as well in the focus of their
372 neuroscientific contributions. Both approaches share the overall goal of using deep probabilistic inference tools to
373 build more interpretable models of neural data. These complementary and concurrent advances will expedite the
374 cycle of building, adjusting and selecting mechanistic models in neuroscience.

375 Finally, a complementary approach to mechanistic modeling is to pursue purely phenomenological models, which
376 are designed to have favorable statistical and computational properties: these data-driven models can be efficiently
377 fit to neural data [18–24, 41, 43] or to implement desired computations [125]. Although tremendously useful for a
378 quantitative characterization of neural dynamics, these models typically have a large number of parameters, which
379 rarely correspond to physically measurable or mechanistically interpretable quantities, and thus it can be challenging
380 to derive mechanistic insights or causal hypotheses from them (but see e.g. [126–128]).

381 Use of summary features

382 When fitting mechanistic models to data, it is common to target summary features to isolate specific behaviors,
383 rather than the full data. For example, the spike shape is known to constrain sodium and potassium conductances
384 [28, 29, 65]. When modeling population dynamics, it is often desirable to achieve realistic firing rates, rate-correlations
385 and response nonlinearities [123, 129], or specified oscillations [7]. In models of decision making, one is often
386 interested in reproducing psychometric functions or reaction-time distributions [130]. Choice of summary features
387 might also be guided by known limitations of either the model or the measurement approach, or necessitated by the
388 fact that published data are only available in summarized form. Several methods have been proposed to automatically
389 construct informative summary features [131–133]. SNPE can be applied to, and might benefit from the use of
390 summary features, but it also makes use of the ability of neural networks to automatically learn informative features
391 in high-dimensional data. Thus, SNPE can also be applied directly to raw data (e.g. using recurrent neural networks
392 [36]), or to high-dimensional summary features which are challenging for ABC approaches (Fig. 2). In all cases, care is
393 needed when interpreting models fit to summary features, as choice of features can influence the results [131–133].

394 Applicability and limitations

395 A key advantage of SNPE is its general applicability: it can be applied whenever one has a simulator that allows to
396 stochastically generate model outputs from specific parameters. Furthermore, it can be applied in a fully ‘black-box’
397 manner’, i.e. does not require access to the internal workings of the simulator, its model equations, likelihoods or
398 gradients. It does not impose any other limitations on the model or the summary features, and in particular does not
399 require them to be differentiable. However, it also has limitations: first, current implementations of SNPE scale well to
400 high-dimensional observations (~ 1000 s dims, also see [37]), but scaling SNPE to even higher-dimensional parameter

401 spaces (>30) is challenging (note that previous approaches were generally limited to $\text{dim} < 10$). Given that the difficulty
402 of estimating full posteriors scales exponentially with dimensionality, this is an inherent challenge for all approaches
403 that aim at full inference (in contrast to just identifying a single, or a few heuristically chosen parameter fits). Second,
404 while it is a long-term goal for these approaches to be made fully automatic, our current implementation still requires
405 choices by the user: as described in Methods, one needs to choose the type of the density estimation network, and
406 specify settings related to network-optimisation, and the number of simulations and inference rounds. These settings
407 depend on the complexity of the relation between summary features and model parameters, and the number of
408 simulations that can be afforded. In the documentation accompanying our code-package, we provide examples and
409 guidance. For small-scale problems, we have found SNPE to be robust to these settings. However, for challenging,
410 high-dimensional applications, SNPE might currently require substantial user interaction. Third, the power of SNPE
411 crucially rests on the ability of deep neural networks to perform density estimation. While deep nets have had ample
412 empirical success, we still have an incomplete understanding of their limitations, in particular in cases where the
413 mapping between data and parameters might not be smooth (e.g. near phase transitions). Fourth, when applying
414 SNPE (or any other model-identification approach), validation of the results is of crucial importance, both to assess
415 the accuracy of the inference procedure, as well as to identify possible limitations of the mechanistic model itself.
416 In the example applications, we used several procedures for assessing the quality of the inferred posteriors. One
417 common ingredient of these approaches is to sample from the inferred model, and search for systematic differences
418 between observed and simulated data, e.g. to perform *posterior predictive checks* [36, 37, 100, 134, 135] (Fig. 2g,
419 Fig. 3f,g, Fig. 4C, and Fig. 5d). There are challenges and opportunities ahead in further scaling and automating
420 simulation-based inference approaches. However, in its current form, SNPE will be a powerful tool for quantitatively
421 evaluating mechanistic hypotheses on neural data, and for designing better models of neural dynamics.

422 Acknowledgments

423 We thank Mahmood S. Hoseini and Michael Stryker for sharing their data for Fig. 2, and Philipp Berens, Sean Bittner,
424 Jan Boelts, John Cunningham, Richard Gao, Scott Linderman, Eve Marder, Iain Murray, George Papamakarios, Astrid
425 Prinz, Auguste Schulz and Srinivas Turaga for discussions and/or comments on the manuscript. This work was
426 supported by the German Research Foundation (DFG) through SFB 1233 ‘Robust Vision’, (276693517), SFB 1089
427 ‘Synaptic Microcircuits’ and SPP 2041 ‘Computational Connectomics’, the German Federal Ministry of Education and
428 Research (BMBF, project ‘ADIMEM’, FKZ 01IS18052 A-D) to JHM, a Sir Henry Dale Fellowship by the Wellcome Trust and
429 the Royal Society (WT100000; WFP and TPV), a Wellcome Trust Senior Research Fellowship (214316/Z/18/Z; TPV), a
430 ERC Consolidator Grant (SYNAPSEEK; WPF and CC), and a UK Research and Innovation, Biotechnology and Biological
431 Sciences Research Council (CC, UKRI-BBSRC BB/N019512/1).

432 **Methods**

433 **Code availability**

434 Code implementing SNPE is available at <http://www.mackelab.org/delfi/>.

435 **Simulation-based inference**

436 To perform Bayesian parameter identification with SNPE, three types of input need to be specified:

- 437 1. A mechanistic model. The model only needs to be specified through a simulator, i.e. that one can generate a
438 simulation result \mathbf{x} for any parameters $\boldsymbol{\theta}$. We do not assume access to the likelihood $p(\mathbf{x}|\boldsymbol{\theta})$ or the equations
439 or internals of the code defining the model, nor do we require the model to be differentiable. This is in
440 contrast to many alternative approaches (including [123]), which require the model to be differentiable and to
441 be implemented in a software code that is amenable to automatic differentiation packages. Finally, SNPE can
442 both deal with inputs \mathbf{x} which resemble ‘raw’ outputs of the model, or summary features calculated from data.
443
- 444 2. Observed data \mathbf{x}_o of the same form as the results \mathbf{x} produced by model simulations.
- 445 3. A prior distribution $p(\boldsymbol{\theta})$ describing the range of possible parameters. $p(\boldsymbol{\theta})$ could consist of upper and lower
446 bounds for each parameter, or a more complex distribution incorporating mechanistic first principles or
knowledge gained from previous inference procedures on other data.

447 For each problem, our goal was to estimate the posterior distribution $p(\boldsymbol{\theta}|\mathbf{x}_o)$. To do this we used SNPE [35–37].

448 Setting up the inference procedure required three design choices:

- 449 1. A network architecture, including number of layers, units per layer, layer type (feedforward or convolutional),
450 activation function and skip connections.
- 451 2. A parametric family of probability densities $q_{\psi}(\boldsymbol{\theta})$ to represent inferred posteriors, to be used as conditional
452 density estimator. We used either a mixture of Gaussians (MoG) or a masked autoregressive flow (MAF) [40]. In
453 the former case, the number of components K must be specified; in the latter the number of MADES (Masked
454 Autoencoder for Distribution Estimation) n_{MADES} . Both choices are able to represent richly structured, and
455 multimodal posterior distributions.
- 456 3. A simulation budget, i.e. number of rounds R and simulations per round N_r .

457 We emphasize that SNPE is highly modular, i.e. that the the inputs (data, the prior over parameter, the mechanistic
458 model), and algorithmic components (network architecture, probability density, optimization approach) can all be
459 modified and chosen independently. This allows neuroscientists to work with models which are designed with
460 mechanistic principles—and not convenience of inference—in mind. Furthermore, it allows SNPE to benefit from
461 advances in more flexible density estimators, more powerful network architectures, or optimization strategies.

462 With the problem and inference settings specified, SNPE adjusts the network weights ϕ based on simulation results,
463 so that $p(\boldsymbol{\theta}|\mathbf{x}) \approx q_{F(\mathbf{x}, \phi)}(\boldsymbol{\theta})$ for any \mathbf{x} . In the first round of SNPE simulation parameters are drawn from the prior $p(\boldsymbol{\theta})$. If
464 a single round of inference is not sufficient, SNPE can be run in multiple rounds, in which samples are drawn from the
465 version of $q_{F(\mathbf{x}_o, \phi)}(\boldsymbol{\theta})$ at the beginning of the round. After the last round, $q_{F(\mathbf{x}_o, \phi)}$ is returned as the inferred posterior on
466 parameters $\boldsymbol{\theta}$ given observed data \mathbf{x}_o . If SNPE is only run for a single round, then the generated samples only depend
467 on the prior, but not on \mathbf{x}_o : in this case, the inference network is applicable to any data (covered by the prior ranges),
468 and can be used for rapid amortized inference.

469 SNPE learns the correct network weights ϕ by minimizing the objective function $\sum_j \mathcal{L}(\boldsymbol{\theta}_j, \mathbf{x}_j)$ where the simulation
470 with parameters $\boldsymbol{\theta}_j$ produced result \mathbf{x}_j . For the first round of SNPE $\mathcal{L}(\boldsymbol{\theta}_j, \mathbf{x}_j) = -\log q_{F(\mathbf{x}, \phi)}$, while in subsequent rounds
471 a different loss function accounts for the fact that simulation parameters were not sampled from the prior. Different
472 choices of the loss function for later rounds result in SNPE-A [35], SNPE-B [36] or SNPE-C algorithm [37]. To optimize
473 the networks, we used ADAM with default settings [136].

474 The details of the algorithm are below:

Algorithm 1: SNPE

Input: simulator with (implicit) density $p(\mathbf{x}|\theta)$, observed data \mathbf{x}_o , prior $p(\theta)$, density family q_ψ , neural network $F(\mathbf{x}, \phi)$, number of rounds R , simulation count for each round N_r

randomly initialize ϕ
 $\tilde{p}_1(\theta) := p(\theta)$
 $N := 0$

for $r = 1$ **to** R **do**

for $i = 1 \dots N_r$ **do**
 | sample $\theta_{N+i} \sim \tilde{p}_r(\theta)$
 | simulate $\mathbf{x}_{N+i} \sim p(\mathbf{x}|\theta_{N+i})$
 $N \leftarrow N + N_r$

train $\phi \leftarrow \arg \min_{\phi} \sum_{j=1}^N \mathcal{L}(\theta_j, \mathbf{x}_j)$

$\tilde{p}_r(\theta) := q_{F(\mathbf{x}_o, \phi)}(\theta)$

return $q_{F(\mathbf{x}_o, \phi)}(\theta)$

476 **Linear-nonlinear encoding models**

We used a Linear-Nonlinear (LN) encoding model (a special case of a generalized linear model, GLM, [18, 20, 41–44]) to simulate the activity of a neuron in response to a univariate time-varying stimulus. Neural activity z_i was subdivided in $T = 100$ bins and, within each bin i , spikes were generated according to a Bernoulli observation model,

$$z_i \sim \text{Bern}(\eta(\mathbf{v}_i^\top \mathbf{f} + \beta)),$$

477 where \mathbf{v}_i is a vector of white noise inputs between time bins $i - 8$ and i , \mathbf{f} a length-9 linear filter, β is the bias, and
 478 $\eta(\cdot) = \exp(\cdot)/(1 + \exp(\cdot))$ is the canonical inverse link function for a Bernoulli GLM. As summary features, we used
 479 the total number of spikes N and the spike-triggered average $\frac{1}{N}\mathbf{Vz}$, where $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_T]$ is the so-called design
 480 matrix of size $9 \times T$. We note that the spike-triggered sum \mathbf{Vz} constitutes sufficient statistics for this GLM, i.e. that
 481 selecting the STA and N together as summary features does not lead to loss of model relevant information over the
 482 full input-output dataset $\{\mathbf{V}, \mathbf{z}\}$. We used a Gaussian prior with zero mean and covariance matrix $\Sigma_\beta = \sigma^2(\mathbf{F}^\top \mathbf{F})^{-1}$,
 483 where \mathbf{F} encourages smoothness by penalizing the second-order differences in the vector of parameters [137].

484 For inference, we used a single round of 10000 simulations, and the posterior was approximated with a Gaussian
 485 distribution ($\theta \in \mathbb{R}^{10}$, $\mathbf{x} \in \mathbb{R}^{10}$). We used a feedforward neural network with two hidden layers of 50 units each. We
 486 used a Polya Gamma Markov Chain Monte Carlo sampling scheme [45] to estimate a reference posterior.

487 In Fig. 2d, we compare the performance of SNPE with two classical ABC algorithms, rejection ABC and Sequential
 488 Monte Carlo ABC as a function of the number of simulations. We report the relative error in Kullback-Leibler divergence,
 489 which is defined as:

$$\frac{D_{\text{KL}}(p_{MCMC}(\theta|\mathbf{x}) \parallel \hat{p}(\theta|\mathbf{x}))}{D_{\text{KL}}(p_{MCMC}(\theta|\mathbf{x}) \parallel p(\theta))}, \quad (1)$$

490 and which ranges between 0 (perfect recovery of the posterior) and 1 (estimated posterior no better than the prior).
 491 Here, $p_{MCMC}(\theta|\mathbf{x})$ is the ground-truth posterior estimated via Markov Chain Monte Carlo sampling, $\hat{p}(\theta|\mathbf{x})$ is the
 492 estimated posterior via SNPE, rejection ABC or Sequential Monte Carlo ABC, and $p(\theta)$ is the prior.

For the spatial receptive field model of a cell in primary visual cortex, we simulated the activity of a neuron
 depending on an image-valued stimulus. Neural activity was subdivided in bins of length $\Delta t = 0.025s$ and within each
 bin i , spikes were generated according to a Poisson observation model,

$$z_i \sim \text{Pois}(\eta(\mathbf{v}_i^\top \mathbf{h} + \beta)),$$

where \mathbf{v}_i is the vectorized white noise stimulus at time bin i , \mathbf{h} a 41×41 linear filter, β is the bias, and $\eta(\cdot) = \exp(\cdot)$ is the canonical inverse link function for a Poisson GLM. The receptive field \mathbf{h} is constrained to be a Gabor filter:

$$h(g_x, g_y) = g \exp\left(-\frac{x'^2 + r^2 y'^2}{2\sigma^2}\right) \cos(2\pi f x' - \phi)$$

$$x' = (g_x - x) \cos \psi - (g_y - y) \sin \psi$$

$$y' = (g_x - x) \sin \psi + (g_y - y) \cos \psi$$

$$\sigma = \frac{\sqrt{2} \log 2}{2\pi f} \frac{2^w + 1}{2^w - 1},$$

493 where (g_x, g_y) is a regular grid of 41×41 positions spanning the 2D image-valued stimulus. The parameters of the
 494 Gabor are gain g , spatial frequency f , aspect-ratio r , width w , phase ϕ (between 0 and π), angle ψ (between 0 and
 495 2π) and location x, y (assumed within the stimulated area, scaled to be between -1 and 1). Bounded parameters
 496 were transformed with a log-, or logit-transform, to yield unconstrained parameters. After applying SNPE, we back-
 497 transformed both the parameters and the estimated posteriors in closed form, as shown in Fig. 2. We did not
 498 transform the bias β .

We used a factorizing Gaussian prior for the vector of transformed Gabor parameters

$$[\log g, \log f, \log r, \log w, l_{0,\pi}(\phi), l_{0,2\pi}(\psi), l_{-1,1}(x), l_{-1,1}(y)],$$

499 where transforms $l_{0,\pi}(X) = \log(X/(2\pi - X))$, $l_{0,2\pi}(X) = \log(X/(\pi - X))$, $l_{-1,1}(X) = \log((X + 1)/(1 - X))$ ensured the
 500 assumed ranges for the Gabor parameters ϕ, ψ, x, y . Our Gaussian prior had zero mean and standard deviations
 501 [0.5, 0.5, 0.5, 0.5, 1.9, 1.78, 1.78, 1.78]. We note that a Gaussian prior on a logit-transformed random variable $\text{logit}X$ with
 502 zero mean and standard deviation around 1.78 is close to a uniform prior over the original variable X . For the bias β ,
 503 we used a Gaussian prior with mean -0.57 and variance 1.63 , which approximately corresponds to an exponential
 504 prior $\exp(\beta) \sim \text{Exp}(\lambda)$ with rate $\lambda = 1$ on the baseline firing rate $\exp(\beta)$ in absence of any stimulus.

505 The ground-truth parameters for the demonstration in Fig. 2 were chosen to give an asymptotic firing rate of 1Hz
 506 for 5 minutes stimulation, resulting in 299 spikes, and a signal-to-noise ratio of -12dB .

507 As summary features, we used the total number of spikes N and the spike-triggered average $\frac{1}{N}\mathbf{Vz}$, where $\mathbf{V} =$
 508 $[v_1, v_2, \dots, v_T]$ is the stimulation video of length $T = 300/\Delta t = 12000$. As for the GLM with a temporal filter, the
 509 spike-triggered sum \mathbf{Vz} constitutes sufficient statistics for this GLM.

510 For inference, we applied SNPE-A with in total 2 rounds: an initial round serves to first roughly identify the
 511 relevant region of parameter space. Here we used a Gaussian distribution to approximate the posterior from 100000
 512 simulations each. A second round then used a mixture of 8 Gaussian components to estimate the exact shape of the
 513 posterior from another 100000 simulations ($\theta \in \mathbb{R}^9, \mathbf{x} \in \mathbb{R}^{1682}$). We used a convolutional network with 5 convolutional
 514 layers with 16 to 32 convolutional filters followed by two fully connected layers with 50 units each. The total number of
 515 spikes N within a simulated experiment was passed as an additional input directly to the fully-connected layers of the
 516 network. Similar to the previous GLM, this model has a tractable likelihood, so we use MCMC to obtain a reference
 517 posterior.

518 We applied this approach to extracellular recordings from primary visual cortex of alert mice obtained using silicon
 519 microelectrodes in response to colored-noise visual stimulation. Experimental methods are described in [51].

520 Comparison with Sequential Monte Carlo (SMC) ABC

521 In order to illustrate the competitive performance of SNPE, we obtained a posterior estimate with a classical ABC
 522 method, Sequential Monte Carlo (SMC) ABC [34, 49]. Likelihood-free inference methods from the ABC family require a
 523 distance function $d(\mathbf{x}_o, \mathbf{x})$ between observed data \mathbf{x}_o and possible simulation outputs \mathbf{x} to characterize dissimilarity
 524 between simulations and data. A common choice is the (scaled) Euclidean distance $d(\mathbf{x}_o, \mathbf{x}) = \|\mathbf{x} - \mathbf{x}_o\|_2$. The Euclidean
 525 distance here was computed over 1681 summary features given by the spike-triggered average (one per pixel) and a
 526 single summary feature given by the ‘spike count’. To ensure that the distance measure was sensitive to differences in
 527 both STA and spike count, we scaled the summary feature ‘spike count’ to account for about 20% of the average total
 528 distance (other values did not yield better results). The other 80% were computed from the remaining 1681 summary
 529 features given by spike-triggered averages. To showcase how this situation is challenging for ABC approaches, we
 530 generated 10000 input-output pairs $(\theta_i, \mathbf{x}_i) \sim p(\mathbf{x}|\theta)p(\theta)$ with the prior and simulator used above, and illustrate the 10

531 STAs and spike counts with closest $d(\mathbf{x}_o, \mathbf{x}_i)$ in Supplementary Fig. 3a. Spike counts were comparable to the observed
 532 data (299 spikes), but STAs are noise-dominated and the 10 ‘closest’ underlying receptive fields (orange contours) show
 533 substantial variability in location and shape of the receptive field. If even the ‘closest’ samples do not show any visible
 534 receptive field, then there is little hope that even an appropriately chosen acceptance threshold will yield a good
 535 approximation to the posterior. These findings were also reflected in the results from SMC-ABC with a total simulation
 536 budget of 10^6 simulations (Fig. 3b). The estimated posterior marginals for ‘bias’ and ‘gain’ parameters show that the
 537 parameters related to the firing rate were constrained by the data \mathbf{x}_o , but marginals of parameters related to shape
 538 and location of the receptive field did not differ from the prior, highlighting that SMC-ABC was not able to identify the
 539 posterior distribution. The low correlations between the ground-truth receptive field and receptive fields sampled
 540 from SMC-ABC posterior further highlight the failure of SMC-ABC to infer the ground-truth posterior (Fig. 3c). Further
 541 comparisons of neural-density estimation approaches with ABC-methods can be found in the studies describing the
 542 underlying machine-learning methodologies [35, 37, 109].

543 **Ion channel models**

We simulated non-inactivating potassium channel currents subject to voltage-clamp protocols as:

$$I_K = \bar{g}_K m(V - E_K),$$

where V is the membrane potential, \bar{g}_K is the density of potassium channels, E_K is the reversal potential of potassium, and m is the gating variable for potassium channel activation. m is modeled according to the first-order kinetic equation

$$\frac{dm}{dt} = \frac{m_\infty(V) - m}{\tau_m(V)},$$

where $m_\infty(V)$ is the steady-state activation, and $\tau_m(V)$ the respective time constant. We used a general formulation of $m_\infty(V)$ and $\tau_m(V)$ [57], where the steady-state activation curve has 2 parameters (slope and offset) and the time constant curve has 6 parameters, amounting to a total of 8 parameters (θ_1 to θ_8):

$$m_\infty(V) = \frac{1}{1 + e^{-\theta_1 V + \theta_2}}$$

$$\tau_m(V) = \frac{\theta_4}{e^{-(\theta_5(V - \theta_3) + \theta_6(V - \theta_3)^2)} + e^{(\theta_7(V - \theta_3) + \theta_8(V - \theta_3)^2)}}.$$

544 Since this model can be used to describe the dynamics of a wide variety of channel models, we refer to it as *Omnimodel*.
 545 We modeled responses of the Omnimodel to a set of five voltage-clamp protocols described in [56]. Current
 546 responses were reduced to 55 summary features (11 per protocol). Summary features were coefficients to basis
 547 functions derived via Principal Components Analysis (PCA) (10 per protocol) plus a linear offset (1 per protocol) found
 548 via least-squares fitting. PCA basis functions were found by simulating responses of the non-inactivating potassium
 549 channel models to the five voltage-clamp protocols and reducing responses to each protocol to 10 dimensions
 550 (explaining 99.9% of the variance).

551 To amortize inference on the model, we specified a wide uniform prior over the parameters: $\theta_1 \in \mathcal{U}(0, 1)$, $\theta_2 \in$
 552 $\mathcal{U}(-10., 10.)$, $\theta_3 \in \mathcal{U}(-120., 120.)$, $\theta_4 \in \mathcal{U}(0., 2000)$, $\theta_5 \in \mathcal{U}(0., 0.5)$, $\theta_6 \in \mathcal{U}(0, 0.05)$, $\theta_7 \in \mathcal{U}(0., 0.5)$, $\theta_8 \in \mathcal{U}(0, 0.05)$.

553 For inference, we trained a shared inference network in a single round of 10^6 simulations generated by sampling
 554 from the prior ($\theta \in \mathbb{R}^8$, $x \in \mathbb{R}^{55}$). The density estimator is a masked autoregressive flow (MAF) [40] with five MADES
 555 with [250,250] hidden units each.

556 We evaluated performance on 350 non-inactivating potassium ion channels selected from IonChannelGenealogy
 557 (ICG) by calculating the correlation coefficient between traces generated by the original model and traces from the
 558 Omnimodel using the posterior mode.

559 **Single-compartment Hodgkin-Huxley neurons**

We simulated a single-compartment Hodgkin-Huxley type neuron with channel kinetics as in [65],

$$C_m \frac{dV}{dt} = g_l(E_l - V) + \bar{g}_{Na} m^3 h (E_{Na} - V) + \bar{g}_K n^4 (E_K - V) + \bar{g}_M p (E_K - V) + I_{inj} + \sigma \eta(t)$$

$$\frac{dq}{dt} = \frac{q_\infty(V) - q}{\tau_q(V)}, q \in \{m, h, n, p\},$$

560 where V is the membrane potential, C_m is the membrane capacitance, g_l is the leak conductance, E_l is the membrane
561 reversal potential, \bar{g}_c is the density of channels of type c (Na^+ , K^+ , M), E_c is the reversal potential of c , (m, h, n, p) are the
562 respective channel gating kinetic variables, and $\sigma\eta(t)$ is the intrinsic neural noise. The right hand side of the voltage
563 dynamics is composed of a leak current, a voltage-dependent Na^+ current, a delayed-rectifier K^+ current, a slow
564 voltage-dependent K^+ current responsible for spike-frequency adaptation, and an injected current I_{inj} . Channel gating
565 variables q have dynamics fully characterized by the neuron membrane potential V , given the respective steady-state
566 $q_\infty(V)$ and time constant $\tau_q(V)$ (details in [65]). Two additional parameters are implicit in the functions $q_\infty(V)$ and
567 $\tau_q(V)$: V_T adjusts the spike threshold through m_∞ , h_∞ , n_∞ , τ_m , τ_h and τ_n ; τ_{\max} scales the time constant of adaptation
568 through $\tau_p(V)$ (details in [65]). We set $E_{\text{Na}} = 53$ mV and $E_K = -107$ mV, similar to the values used for simulations in
569 Allen Cell Types Database (<http://help.brain-map.org/download/attachments/8323525/BiophysModelPeri.pdf>).

570 We applied SNPE to infer the posterior over 8 parameters ($\bar{g}_{\text{Na}}, \bar{g}_K, g_l, \bar{g}_M, \tau_{\max}, V_T, \sigma, E_l$), given 7 voltage features
571 (number of spikes, mean resting potential, standard deviation of the resting potential, and the first 4 voltage moments,
572 mean, standard deviation, skewness and kurtosis).

The prior distribution over the parameters was uniform,

$$\boldsymbol{\theta} \sim \mathcal{U}(p_{\text{low}}, p_{\text{high}}),$$

573 where $p_{\text{low}} = [0.5, 10^{-4}, 10^{-4}, 10^{-4}, 50, 40, 10^{-4}, 35]$ and $p_{\text{high}} = [80, 15, 0.6, 0.6, 3000, 90, 0.15, 100]$. These ranges are
574 similar to the ones obtained in [65].

575 For inference in simulated data, we used a single round of 100000 simulations ($\boldsymbol{\theta} \in \mathbb{R}^8, \mathbf{x} \in \mathbb{R}^{11}$). The density
576 estimator was a masked autoregressive flow (MAF) [40] with five MADEs with [50,50] hidden units each.

577 For the inference on in vitro recordings from mouse cortex (Allen Cell Types Database, <https://celltypes.brain-map.org/data>), we selected 8 recordings corresponding to spiny neurons with at least 10 spikes during the current-
578 clamp stimulation. The respective cell identities and sweeps are: (518290966,57), (509881736,39), (566517779,46),
579 (567399060,38), (569469018,44), (532571720,42), (555060623,34), (534524026,29). For each recording, SNPE-B was run
580 for 2 rounds with 125000 Hodgkin-Huxley simulations each, and the posterior was approximated by a mixture of two
581 Gaussians. In this case, the density estimator was composed of two fully connected layers of 100 units each.

583 Comparison with genetic algorithm

We compared SNPE posterior with a state-of-the-art genetic algorithm (Indicator Based Evolutionary Algorithm IBEA, [70, 71] from the BluePyOpt package [31]), in the context of the Hodgkin-Huxley model with 8 parameters and 7 features (Supplementary Fig. 9). For each Hodgkin-Huxley model simulation i and summary feature j , we used the following objective score:

$$\epsilon_{ij} = \left| \frac{x_{ij} - x_{oj}}{\sigma_j} \right|, \quad j = 1, \dots, 7,$$

584 where x_{ij} is the value of summary feature j for simulation i , x_{oj} is the observed summary feature j , and σ_j is the
585 standard deviation of the summary feature j computed across 1000 previously simulated datasets. IBEA outputs the
586 hall-of-fame, which corresponds to the 10 parameter sets with the lowest sum of objectives $\sum_j \epsilon_{ij}$. We ran IBEA with
587 100 generations and an offspring size of 1000 individuals, corresponding to a total of 100000 simulations.

588 Circuit model of the crustacean stomatogastric ganglion

589 We used extracellular nerve recordings made from the stomatogastric motor neurons that principally comprise the
590 triphasic pyloric rhythm in the crab *Cancer borealis* [74]. The preparations were decentralized, i.e. the axons of the
591 descending modulatory inputs were severed. The data was recorded at a temperature of 11 °C. See [74] for full
592 experimental details.

We simulated the circuit model of the crustacean stomatogastric ganglion by adapting a model described in [7]. The model is composed of three single-compartment neurons, AB/PD, LP, and PD, where the electrically coupled AB and PD neurons are modeled as a single neuron. Each of the model neurons contains 8 currents, a Na^+ current I_{Na} , a fast and a slow transient Ca^{2+} current I_{CaT} and I_{CaS} , a transient K^+ current I_A , a Ca^{2+} -dependent K^+ current I_{KCa} , a delayed rectifier K^+ current I_{Kd} , a hyperpolarization-activated inward current I_H , and a leak current I_{leak} . In addition, the model contains 7 synapses. As in [7], these synapses were simulated using a standard model of synaptic dynamics [138]. The synaptic input current into the neurons is given by $I_s = g_{ss}(V_{\text{post}} - E_s)$, where g_s is the maximal synapse

conductance, V_{post} the membrane potential of the postsynaptic neuron, and E_s the reversal potential of the synapse. The evolution of the activation variable s is given by

$$\frac{ds}{dt} = \frac{\bar{s}(V_{\text{pre}}) - s}{\tau_s}$$

with

$$\bar{s}(V_{\text{pre}}) = \frac{1}{1 + \exp((V_{\text{th}} - V_{\text{pre}})/\delta)} \quad \text{and} \quad \tau_s = \frac{1 - \bar{s}(V_{\text{pre}})}{k_-}.$$

593 Here, V_{pre} is the membrane potential of the presynaptic neuron, V_{th} is the half-activation voltage of the synapse, δ sets
594 the slope of the activation curve, and k_- is the rate constant for transmitter-receptor dissociation rate.

595 As in [7], two types of synapses were modeled since AB, LP, and PY are glutamatergic neurons whereas PD is
596 cholinergic. We set $E_s = -70$ mV and $k_- = 1/40$ ms for all glutamatergic synapses and $E_s = -80$ mV and $k_- = 1/100$
597 ms for all cholinergic synapses. For both synapse types, we set $V_{\text{th}} = -35$ mV and $\delta = 5$ mV.

598 For each set of membrane and synaptic conductances, we numerically simulated the rhythm for 10 seconds with
599 a step size of 0.025 ms. To make the model stochastic, at each time step, we added Gaussian noise with a standard
600 deviation of 0.001 mV to the input of each neuron.

601 We applied SNPE to infer the posterior over 24 membrane parameters and 7 synaptic parameters, i.e. 31 pa-
602 rameters in total. The 7 synaptic parameters were the maximal conductances g_s of all synapses in the circuit,
603 each of which is varied uniformly in logarithmic domain from 0.01 nS to 1000 nS, with an exception of the synapse
604 from AB to LP, which is varied uniformly in logarithmic domain from 0.01 nS to 10000 nS. The membrane pa-
605 rameters were the maximal membrane conductances for each of the neurons. The membrane conductances were
606 varied over an extended range of previously reported values [7], which led us to the uniform prior bounds $p_{\text{low}} =$
607 $[0, 0, 0, 0, 0, 25, 0, 0] \text{ mS cm}^{-2}$ and $p_{\text{high}} = [500, 7.5, 8, 60, 15, 150, 0.2, 0.01] \text{ mS cm}^{-2}$ for the maximal membrane con-
608 ductances of the AB neuron, $p_{\text{low}} = [0, 0, 2, 10, 0, 0, 0, 0.01] \text{ mS cm}^{-2}$ and $p_{\text{high}} = [200, 2.5, 12, 60, 10, 125, 0.06, 0.04] \text{ mS cm}^{-2}$
609 for the maximal membrane conductances of the LP neuron, and $p_{\text{low}} = [0, 0, 0, 30, 0, 50, 0, 0] \text{ mS cm}^{-2}$ and $p_{\text{high}} =$
610 $[600, 12.5, 4, 60, 5, 150, 0.06, 0.04] \text{ mS cm}^{-2}$ for the maximal membrane conductances of the PY neuron. The order of
611 the membrane currents was: [Na, CaT, CaS, A, KCa, Kd, H, leak].

612 We used the 15 summary features proposed by [7], and extended them by 3 additional features. The features
613 proposed by [7] are 15 salient features of the pyloric rhythm, namely: cycle period T (s), AB/PD burst duration d_{AB}^b (s),
614 LP burst duration d_{LP}^b (s), PY burst duration d_{PY}^b (s), gap AB/PD end to LP start $\Delta t_{\text{AB-LP}}^{\text{es}}$ (s), gap LP end to PY start $\Delta t_{\text{LP-PY}}^{\text{es}}$
615 (s), delay AB/PD start to LP start $\Delta t_{\text{AB-LP}}^{\text{ss}}$ (s), delay LP start to PY start $\Delta t_{\text{LP-PY}}^{\text{ss}}$ (s), AB/PD duty cycle d_{AB} , LP duty cycle d_{LP} ,
616 PY duty cycle d_{PY} , phase gap AB/PD end to LP start $\Delta\phi_{\text{AB-LP}}$, phase gap LP end to PY start $\Delta\phi_{\text{LP-PY}}$, LP start phase ϕ_{LP} ,
617 and PY start phase ϕ_{PY} . Note that several of these values are only defined if each neuron produces rhythmic bursting
618 behavior. In addition, for each of the three neurons, we used one feature that describes the maximal duration of its
619 voltage being above -30 mV. We did this as we observed plateaus at around -10 mV during the onset of bursts, and
620 wanted to distinguish such activity traces from others. If the maximal duration was below 5 ms, we set this feature to 5
621 ms. To extract the summary features from the observed experimental data, we first found spikes by searching for
622 local maxima above a hand-picked voltage threshold, and then extracted the 15 above described features. We set the
623 additional 3 features to 5 ms.

624 We used SNPE to infer the posterior distribution over the 18 summary features from experimental data. For
625 inference, we used a single round with 18.5 million samples, out of which 174,000 samples contain bursts in all
626 neurons. We therefore used these 174,000 samples with well defined summary features for training the inference
627 network ($\theta \in \mathbb{R}^{31}, \mathbf{x} \in \mathbb{R}^{18}$). The density estimator was a masked autoregressive flow (MAF) [40] with five MADEs with
628 [200,400] hidden units each. The synaptic conductances were transformed into logarithmic space before training and
629 for the entire analysis.

630 Previous approaches for fitting the STG circuit [7] first fit individual neuron features and reduce the number of
631 possible neuron models [25], and then fit the whole circuit model. While powerful, this approach both requires the
632 availability of single-neuron data, and cannot give access to potential compensation mechanisms between single-
633 neuron and synaptic parameters. Unlike [7], we apply SNPE to directly identify the full 31 dimensional parameter space
634 without requiring experimental measurements of each individual neuron in the circuit. Despite the high-dimensional
635 parameter space, SNPE can identify the posterior distribution using 18 million samples, whereas a direct application

636 of a full-grid method would require $4.65 \cdot 10^{21}$ samples to fill the 31 dimensional parameter space on a grid with five
637 values per dimension.

638 Finding paths in the posterior

639 In order to find directions of robust network output, we searched for a path of high posterior probability. First, as in
640 [7], we aimed to find 2 similar model outputs with disparate parameters. To do so, we sampled from the posterior and
641 searched for 2 parameter sets whose summary features were within 0.1 standard deviations of all 174,000 samples
642 from the observed experimental data, but that had strongly disparate parameters from each other. In the following,
643 we denote the obtained parameter sets by θ_s and θ_g .

644 Second, in order to identify whether network output can be maintained along a continuous path between these 2
645 samples, we searched for a connection in parameter space lying in regions of high posterior probability. To do so, we
646 considered the connection between the samples as a path and minimize the following path integral:

$$\mathcal{L}(\gamma) = \int_0^1 -\log(p_{\theta|x}(\gamma(s))|\mathbf{x}_o) \|\dot{\gamma}(s)\| ds. \quad (2)$$

To minimize this term, we parameterized the path $\gamma(s)$ using sinusoidal basis-functions with coefficients $\alpha_{n,k}$:

$$\gamma(s) = \begin{bmatrix} \sum_{k=1}^K \alpha_{1,k} \cdot \sin(\pi ks) \\ \vdots \\ \sum_{k=1}^K \alpha_{N,k} \cdot \sin(\pi ks) \end{bmatrix} + \begin{bmatrix} \sum_{k=K+1}^{2K} \alpha_{1,k} \cdot \sin^2(\pi ks) \\ \vdots \\ \sum_{k=K+1}^{2K} \alpha_{N,k} \cdot \sin^2(\pi ks) \end{bmatrix} + (1-s) \cdot \theta_s + s\theta_g$$

These basis functions are defined such that, for any coefficients $\alpha_{n,k}$, the starting and end points of the path are exactly
the two parameter sets defined above:

$$\gamma(0) = \theta_s \quad \gamma(1) = \theta_g$$

647 With this formulation, we have framed the problem of finding the path as an unconstrained optimization problem over
648 the parameters $\alpha_{n,k}$. We can therefore minimize the path integral L using gradient descent over $\alpha_{n,k}$. For numerical
649 simulations, we approximated the integral in equation 2 as a sum over 80 points along the path and use 2 basis
650 functions for each of the 31 dimensions, i.e. $K = 2$.

In order to demonstrate the sensitivity of the pyloric network, we aimed to find a path along which the circuit
output quickly breaks down. For this, we picked a starting point along the high-probability path and then minimize
the posterior probability. In addition, we enforced that the orthogonal path lies within an orthogonal disk to the
high-probability path, leading to the following constrained optimization problem:

$$\min_{\theta} \log(p(\theta|x)) \quad \text{s.t.} \quad n^T \Delta\theta = 0$$

where n is the tangent vector along the path of high probability. This optimization problem can be solved using the
gradient projection method [139]:

$$\Delta\theta = -\frac{P(\nabla \log(p(\theta|x)))}{\sqrt{(\nabla \log(p(\theta|x)))^T P(\nabla \log(p(\theta|x)))}}$$

651 with projection matrix $P = \mathbb{1} - \frac{1}{n^T n} nn^T$ and $\mathbb{1}$ indicating the identity matrix. Each gradient update is a step along
652 the orthogonal path. We let the optimization run until the distance along the path is 1/27 of the distance along the
653 high-probability path.

654 Identifying conditional correlations

655 In order to investigate compensation mechanisms in the STG, we compared marginal and conditional correlations.
656 For the marginal correlation matrix in Fig. 6b, we calculated the Pearson correlation coefficient based on 1.26 million
657 samples from the posterior distribution $p(\theta|x)$. To find the 2-dimensional conditional distribution for any pair of
658 parameters, we fixed all other parameters to values taken from an arbitrary posterior sample, and varied the remaining
659 2 on an evenly spaced grid with 50 points along each dimension, covering the entire prior space. We evaluated the
660 posterior distribution at every value on this grid. We then calculated the conditional correlation as the Pearson
661 correlation coefficient over this distribution. For the 1-dimensional conditional distribution, we varied only 1 parameter
662 and kept all others fixed. Lastly, in Fig. 6d, we sampled 500 parameter sets from the posterior, computed the respective
663 conditional posteriors and conditional correlation matrices, and took the average over the conditional correlation
664 matrices.

665 References

- 666 [1] A. V. Herz, T. Gollisch, C. K. Machens, and D. Jaeger. Modeling single-neuron dynamics and computations: a balance of detail
667 and abstraction. *Science*, 314(5796):80–85, 2006.
- 668 [2] W. Gerstner, H. Sprekeler, and G. Deco. Theory and simulation in neuroscience. *Science*, 338(6103):60–65, 2012.
- 669 [3] T. O’Leary, A. C. Sutton, and E. Marder. Computational models in the age of large datasets. *Current Opinion in Neurobiology*, 32:
670 87–94, 2015.
- 671 [4] R. E. Baker, J.-M. Pena, J. Jayamohan, and A. Jérusalem. Mechanistic models versus machine learning, a fight worth fighting for
672 the biological community? *Biology Letters*, 14(5), 2018.
- 673 [5] A. L. Hodgkin and A. F. Huxley. A quantitative description of membrane current and its application to conduction and excitation
674 in nerve. *The Journal of Physiology*, 117(4):500–544, 1952.
- 675 [6] C. van Vreeswijk and H. Sompolinsky. Chaos in neuronal networks with balanced excitatory and inhibitory activity. *Science*, 274
676 (5293), 1996.
- 677 [7] A. A. Prinz, D. Bucher, and E. Marder. Similar network activity from disparate circuit parameters. *Nature Neuroscience*, 7(12):
678 1345, 2004.
- 679 [8] T. P. Vogels, K. Rajan, and L. F. Abbott. Neural network dynamics. *Annual Review of Neuroscience*, 28:357–376, 2005.
- 680 [9] T. C. Potjans and M. Diesmann. The cell-type specific cortical microcircuit: relating structure and activity in a full-scale spiking
681 network model. *Cerebral Cortex*, 24(3):785–806, 2012.
- 682 [10] A. Litwin-Kumar and B. Doiron. Slow dynamics and high variability in balanced cortical networks with clustered connections.
683 *Nature Neuroscience*, 15(11):1498, 2012.
- 684 [11] O. Sporns. Contributions and challenges for network models in cognitive neuroscience. *Nature Neuroscience*, 17(5):652, 2014.
- 685 [12] D. S. Bassett, P. Zurn, and J. I. Gold. On the nature and use of models in network neuroscience. *Nature Reviews Neuroscience*, 19
686 (9):566, 2018.
- 687 [13] J. I. Gold and M. N. Shadlen. The neural basis of decision making. *Annual Review of Neuroscience*, 30, 2007.
- 688 [14] X.-J. Wang. Decision making in recurrent neuronal circuits. *Neuron*, 60(2):215–234, 2008.
- 689 [15] R. N. Gutenkunst, J. J. Waterfall, F. P. Casey, K. S. Brown, C. R. Myers, and J. P. Sethna. Universally sloppy parameter sensitivities
690 in systems biology models. *PLoS Computational Biology*, 3(10):e189, 2007.
- 691 [16] P. Achard and E. De Schutter. Complex parameter landscape for a complex neuron model. *PLoS Computational Biology*, 2(7):
692 e94, 2006.
- 693 [17] L. M. Alonso and E. Marder. Visualization of currents in neural models with similar behavior and different conductance
694 densities. *eLife*, 8:e42722, 2019.
- 695 [18] W. Truccolo, U. T. Eden, M. R. Fellows, J. P. Donoghue, and E. N. Brown. A point process framework for relating neural spiking
696 activity to spiking history, neural ensemble, and extrinsic covariate effects. *Journal of Neurophysiology*, 93(2):1074–1089, 2005.
- 697 [19] E. Schneidman, M. J. Berry II, R. Segev, and W. Bialek. Weak pairwise correlations imply strongly correlated network states in a
698 neural population. *Nature*, 440(7087):1007, 2006.
- 699 [20] J. W. Pillow, J. Shlens, L. Paninski, A. Sher, A. M. Litke, E. J. Chichilnisky, and E. P. Simoncelli. Spatio-temporal correlations and
700 visual signalling in a complete neuronal population. *Nature*, 454(7207), 2008.
- 701 [21] B. M. Yu, J. P. Cunningham, G. Santhanam, S. I. Ryu, K. V. Shenoy, and M. Sahani. Gaussian-process factor analysis for
702 low-dimensional single-trial analysis of neural population activity. *Journal of Neurophysiology*, 102(1):614–35, Jul 2009.
703 doi:10.1152/jn.90941.2008.
- 704 [22] J. H. Macke, L. Buesing, J. P. Cunningham, B. M. Yu, K. V. Shenoy, and M. Sahani. Empirical models of spiking in neural
705 populations. In *Advances in Neural Information Processing Systems*, pages 1350–1358, 2011.
- 706 [23] J. P. Cunningham and B. M. Yu. Dimensionality reduction for large-scale neural recordings. *Nature Neuroscience*, 17(11):1500,
707 2014.
- 708 [24] C. Pandarinath, D. J. O’Shea, J. Collins, R. Jozefowicz, S. D. Stavisky, J. C. Kao, E. M. Trautmann, M. T. Kaufman, S. I. Ryu, L. R.
709 Hochberg, et al. Inferring single-trial neural population dynamics using sequential auto-encoders. *Nature Methods*, page 1,
710 2018.

- 711 [25] A. A. Prinz, C. P. Billimoria, and E. Marder. Alternative to hand-tuning conductance-based models: construction and analysis of
712 databases of model neurons. *Journal of Neurophysiology*, 90(6):3998–4015, 2003.
- 713 [26] C. Tomm, M. Avermann, T. Vogels, W. Gerstner, and C. Petersen. The influence of structure on the response properties of
714 biologically plausible neural network models. *BMC neuroscience*, 12(1):P30, 2011.
- 715 [27] C. Stringer, M. Pachitariu, N. A. Steinmetz, M. Okun, P. Bartho, K. D. Harris, M. Sahani, and N. A. Lesica. Inhibitory control of
716 correlated intrinsic variability in cortical networks. *eLife*, 5, 2016.
- 717 [28] S. Druckmann, Y. Banitt, A. A. Gidon, F. Schürmann, H. Markram, and I. Segev. A novel multiple objective optimization framework
718 for constraining conductance-based neuron models by experimental data. *Frontiers in Neuroscience*, 1:1, 2007.
- 719 [29] E. Hay, S. Hill, F. Schürmann, H. Markram, and I. Segev. Models of neocortical layer 5b pyramidal cells capturing a wide range of
720 dendritic and perisomatic active properties. *PLoS Computational Biology*, 7(7), 2011.
- 721 [30] C. Rossant, D. F. M. Goodman, B. Fontaine, J. Platkiewicz, A. K. Magnusson, and R. Brette. Fitting neuron models to spike trains.
722 *Frontiers in Neuroscience*, 5:9, 2011.
- 723 [31] W. Van Geit, M. Gevaert, G. Chindemi, C. Rössert, J. Courcol, E. B. Muller, F. Schürmann, I. Segev, and H. Markram. Bluepyopt:
724 Leveraging open source software and cloud infrastructure to optimise model parameters in neuroscience. *Frontiers in
725 Neuroinformatics*, 10:17, 2016.
- 726 [32] M. Beaumont, W. Zhang, and D. J. Balding. Approximate bayesian computation in population genetics. *Genetics*, 162(4), 2002.
- 727 [33] P. Marjoram, J. Molitor, V. Plagnol, and S. Tavaré. Markov chain monte carlo without likelihoods. *Proceedings of the National
728 Academy of Sciences*, 100(26), 2003.
- 729 [34] S. A. Sisson, Y. Fan, and M. M. Tanaka. Sequential monte carlo without likelihoods. *Proceedings of the National Academy of
730 Sciences*, 104(6):1760–1765, 2007.
- 731 [35] G. Papamakarios and I. Murray. Fast ϵ -free inference of simulation models with bayesian conditional density estimation. In
732 *Advances in Neural Information Processing Systems*, pages 1028–1036, 2016.
- 733 [36] J.-M. Lueckmann, P. J. Goncalves, G. Bassetto, K. Öcal, M. Nonnenmacher, and J. H. Macke. Flexible statistical inference for
734 mechanistic models of neural dynamics. In *Advances in Neural Information Processing Systems*, pages 1289–1299, 2017.
- 735 [37] D. Greenberg, M. Nonnenmacher, and J. Macke. Automatic posterior transformation for likelihood-free inference. In *International
736 Conference on Machine Learning*, pages 2404–2414, 2019.
- 737 [38] K. Cranmer, J. Brehmer, and G. Louppe. The frontier of simulation-based inference. *arXiv preprint arXiv:1911.01429*, 2019.
- 738 [39] D. J. Rezende and S. Mohamed. Variational inference with normalizing flows. In *Proceedings of the 32nd International Conference
739 on International Conference on Machine Learning-Volume 37*, pages 1530–1538. JMLR. org, 2015.
- 740 [40] G. Papamakarios, T. Pavlakou, and I. Murray. Masked autoregressive flow for density estimation. In *Advances in Neural
741 Information Processing Systems*, pages 2338–2347, 2017.
- 742 [41] E. N. Brown, L. M. Frank, D. Tang, M. C. Quirk, and M. A. Wilson. A statistical paradigm for neural spike train decoding applied to
743 position prediction from ensemble firing patterns of rat hippocampal place cells. *Journal of Neuroscience*, 18(18):7411–7425,
744 1998.
- 745 [42] L. Paninski. Maximum likelihood estimation of cascade point-process neural encoding models. *Network: Computation in Neural
746 Systems*, 15(4):243–262, 2004.
- 747 [43] J. Pillow. Likelihood-based approaches to modeling the neural code. *Bayesian Brain: Probabilistic Approaches to Neural Coding*,
748 pages 53–70, 2007.
- 749 [44] S. Gerwinn, J. H. Macke, and M. Bethge. Bayesian inference for generalized linear models for spiking neurons. *Frontiers in
750 Computational Neuroscience*, 4:12, 2010.
- 751 [45] N. G. Polson, J. G. Scott, and J. Windle. Bayesian inference for logistic models using pólya–gamma latent variables. *Journal of
752 the American Statistical Association*, 108(504):1339–1349, 2013.
- 753 [46] J. W. Pillow and J. Scott. Fully bayesian inference for neural models with negative-binomial spiking. In *Advances in Neural
754 Information Processing Systems*, pages 1898–1906, 2012.
- 755 [47] J. W. Pillow, L. Paninski, V. J. Uzzell, E. P. Simoncelli, and E. Chichilnisky. Prediction and decoding of retinal ganglion cell responses
756 with a probabilistic spiking model. *Journal of Neuroscience*, 25(47):11003–11013, 2005.

- 757 [48] E. Chichilnisky. A simple white noise analysis of neuronal light responses. *Network: Computation in Neural Systems*, 12(2):
758 199–213, 2001.
- 759 [49] M. A. Beaumont, J. Cornuet, J. Marin, and C. P. Robert. Adaptive approximate bayesian computation. *Biometrika*, 2009.
- 760 [50] C. M. Niell and M. P. Stryker. Highly selective receptive fields in mouse visual cortex. *Journal of Neuroscience*, 28(30):7520–7536,
761 2008. ISSN 0270-6474. doi:10.1523/JNEUROSCI.0623-08.2008.
- 762 [51] L. Dyballa, M. S. Hoseini, M. C. Dadarlat, S. W. Zucker, and M. P. Stryker. Flow stimuli reveal ecologically appropriate responses
763 in mouse visual cortex. *Proceedings of the National Academy of Sciences*, 115(44):11304–11309, 2018. ISSN 0027-8424.
764 doi:10.1073/pnas.1811265115.
- 765 [52] J. P. Jones and L. A. Palmer. An evaluation of the two-dimensional gabor filter model of simple receptive fields in cat striate
766 cortex. *Journal of Neurophysiology*, 58(6):1233–1258, 1987.
- 767 [53] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in
768 Neural Information Processing Systems*, pages 1097–1105, 2012.
- 769 [54] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference
770 on Learning Representations*, 2015.
- 771 [55] C. M. Bishop. Mixture density networks. *Technical Report. Aston University, Birmingham*, 1994.
- 772 [56] W. F. Podlaski, A. Seeholzer, L. N. Groschner, G. Miesenböck, R. Ranjan, and T. P. Vogels. Mapping the function of neuronal ion
773 channels in model and experiment. *eLife*, 6:e22152, March 2017. ISSN 2050-084X. doi:10.7554/eLife.22152.
- 774 [57] A. Destexhe and J. R. Huguenard. Nonlinear Thermodynamic Models of Voltage-Dependent Currents. *Journal of Computational
775 Neuroscience*, 9(3):259–270, November 2000. ISSN 1573-6873. doi:10.1023/A:1026535704537.
- 776 [58] J. Dunlop, M. Bowlby, R. Peri, D. Vasilyev, and R. Arias. High-throughput electrophysiology: an emerging paradigm for
777 ion-channel screening and physiology. *Nature Reviews Drug Discovery*, 7(4):358, 2008.
- 778 [59] H.-J. Suk, E. S. Boyden, and I. van Welie. Advances in the automation of whole-cell patch clamp technology. *Journal of
779 Neuroscience Methods*, 326:108357, 2019. ISSN 0165-0270. doi:<https://doi.org/10.1016/j.jneumeth.2019.108357>.
- 780 [60] R. Ranjan, E. Logette, M. Marani, M. Herzog, V. Tache, and H. Markram. A kinetic map of the homomeric voltage-gated
781 potassium channel (kv) family. *Frontiers in Cellular Neuroscience*, 13:358, 2019.
- 782 [61] A. Speiser, J. Yan, E. W. Archer, L. Buesing, S. C. Turaga, and J. H. Macke. Fast amortized inference of neural activity from calcium
783 imaging data with variational autoencoders. In *Advances in Neural Information Processing Systems*, pages 4024–4034, 2017.
- 784 [62] S. Webb, A. Golinski, R. Zinkov, S. Narayanaswamy, T. Rainforth, Y. W. Teh, and F. Wood. Faithful inversion of generative models
785 for effective amortized inference. In *Advances in Neural Information Processing Systems*, pages 3070–3080, 2018.
- 786 [63] T. S. McTavish, M. Migliore, G. M. Shepherd, and M. L. Hines. Mitral cell spike synchrony modulated by dendrodendritic synapse
787 location. *Frontiers in computational neuroscience*, 6:3, 2012.
- 788 [64] Q. J. M. Huys, M. B. Ahrens, and L. Paninski. Efficient estimation of detailed single-neuron models. *Journal of Neurophysiology*,
789 96(2), 2006.
- 790 [65] M. Pospischil, M. Toledo-Rodriguez, C. Monier, Z. Piwkowska, T. Bal, Y. Frégnac, H. Markram, and A. Destexhe. Minimal
791 hodgkin-huxley type models for different classes of cortical and thalamic neurons. *Biological Cybernetics*, 99(4-5), 2008.
- 792 [66] C. D. Meliza, M. Kostuk, H. Huang, A. Nogaret, D. Margoliash, and H. D. Abarbanel. Estimating parameters and predicting
793 membrane voltages with conductance-based neuron models. *Biological Cybernetics*, 108(4):495–516, 2014.
- 794 [67] R. Ben-Shalom, J. Balewski, A. Siththaranjan, V. Baratham, H. Kyoung, K. G. Kim, K. J. Bender, and K. E. Bouchard. Inferring
795 neuronal ionic conductances from membrane potentials using cnns. *bioRxiv*, 2019. doi:10.1101/727974.
- 796 [68] A. C. Daly, D. J. Gavaghan, C. Holmes, and J. Cooper. Hodgkin–huxley revisited: reparametrization and identifiability analysis of
797 the classic action potential model with approximate bayesian methods. *Royal Society Open Science*, 2(12):150499, 2015.
- 798 [69] N. W. Gouwens, J. Berg, D. Feng, S. A. Sorensen, H. Zeng, M. J. Hawrylycz, C. Koch, and A. Arkhipov. Systematic generation of
799 biophysically detailed models for diverse cortical neuron types. *Nature Communications*, 9(1):710, 2018.
- 800 [70] S. Bleuler, M. Laumanns, L. Thiele, and E. Zitzler. Pisa—a platform and programming language independent interface for
801 search algorithms. In *International Conference on Evolutionary Multi-Criterion Optimization*, pages 494–508. Springer, 2003.

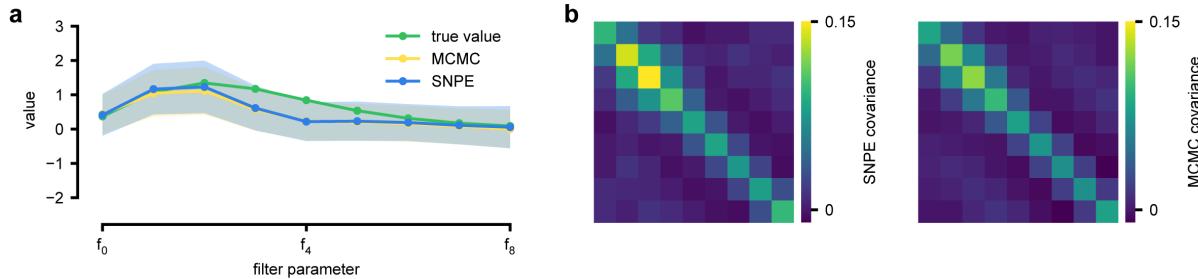
- 802 [71] E. Zitzler and S. Künzli. Indicator-based selection in multiobjective search. In *International conference on parallel problem solving from nature*, pages 832–842. Springer, 2004.
- 803
- 804 [72] Allen Institute for Brain Science. Allen cell types database. <http://celltypes.brain-map.org/>, 2016.
- 805 [73] C. Teeter, R. Iyer, V. Menon, N. Gouwens, D. Feng, J. Berg, A. Szafer, N. Cain, H. Zeng, M. Hawrylycz, et al. Generalized leaky
806 integrate-and-fire models classify multiple neuron types. *Nature Communications*, 9(1):709, 2018.
- 807 [74] S. A. Haddad and E. Marder. Circuit robustness to temperature perturbation is altered by neuromodulators. *Neuron*, 100(3):
808 609–623, 2018.
- 809 [75] A. L. Taylor, T. J. Hickey, A. A. Prinz, and E. Marder. Structure and visualization of high-dimensional conductance spaces. *Journal
810 of Neurophysiology*, 96(2):891–905, 2006.
- 811 [76] E. Marder and J.-M. Goaillard. Variability, compensation and homeostasis in neuron and network function. *Nature Reviews
812 Neuroscience*, 7(7):563, 2006.
- 813 [77] G. J. Gutierrez, T. O’Leary, and E. Marder. Multiple mechanisms switch an electrically coupled, synaptically inhibited neuron
814 between competing rhythmic oscillators. *Neuron*, 77(5):845–858, 2013.
- 815 [78] D. Fisher, I. Olasagasti, D. W. Tank, E. R. Aksay, and M. S. Goldman. A modeling framework for deriving the structural and
816 functional architecture of a short-term memory microcircuit. *Neuron*, 79(5):987–1000, 2013.
- 817 [79] E. Marder, M. L. Goeritz, and A. G. Otopalik. Robust circuit rhythms in small circuits arise from variable circuit components and
818 mechanisms. *Current Opinion in Neurobiology*, 31:156–163, 2015.
- 819 [80] T. O’Leary, A. H. Williams, A. Franci, and E. Marder. Cell types, network homeostasis, and pathological compensation from a
820 biologically plausible ion channel expression model. *Neuron*, 82(4):809–821, 2014.
- 821 [81] T. O’Leary and E. Marder. Temperature-robust neural function from activity-dependent ion channel regulation. *Current Biology*,
822 26(21):2935–2941, 2016.
- 823 [82] M. S. Goldman, J. Golowasch, E. Marder, and L. Abbott. Global structure, robustness, and modulation of neuronal models.
824 *Journal of Neuroscience*, 21(14):5229–5238, 2001.
- 825 [83] B. B. Machta, R. Chachra, M. K. Transtrum, and J. P. Sethna. Parameter space compression underlies emergent theories and
826 predictive models. *Science*, 342(6158):604–607, 2013.
- 827 [84] J. N. MacLean, Y. Zhang, M. L. Goeritz, R. Casey, R. Oliva, J. Guckenheimer, and R. M. Harris-Warrick. Activity-independent
828 coregulation of ia and ih in rhythmically active neurons. *Journal of Neurophysiology*, 94(5):3601–3617, 2005.
- 829 [85] R. Grashow, T. Brookings, and E. Marder. Compensation for variable intrinsic neuronal excitability by circuit-synaptic interactions.
830 *Journal of Neuroscience*, 30(27):9145–9156, 2010.
- 831 [86] E. Marder. Variability, compensation, and modulation in neurons and circuits. *Proceedings of the National Academy of Sciences*,
832 108(Supplement 3):15542–15548, 2011.
- 833 [87] E. Marder and A. L. Taylor. Multiple models to capture the variability in biological neurons and networks. *Nature Neuroscience*,
834 14(2):133, 2011.
- 835 [88] T. O’Leary. Homeostasis, failure of homeostasis and degenerate ion channel regulation. *Current Opinion in Physiology*, 2:
836 129–138, 2018.
- 837 [89] A. L. Taylor, J.-M. Goaillard, and E. Marder. How multiple conductances determine electrophysiological properties in a
838 multicompartment model. *Journal of Neuroscience*, 29(17):5573–5586, 2009.
- 839 [90] J. Golowasch, M. S. Goldman, L. Abbott, and E. Marder. Failure of averaging in the construction of a conductance-based neuron
840 model. *Journal of neurophysiology*, 87(2):1129–1131, 2002.
- 841 [91] H. Kitano. Biological robustness. *Nature Reviews Genetics*, 5(11):826–837, 2004.
- 842 [92] J. N. MacLean, Y. Zhang, B. R. Johnson, and R. M. Harris-Warrick. Activity-independent homeostasis in rhythmically active
843 neurons. *Neuron*, 37(1):109–120, 2003.
- 844 [93] A. V. M. Herz, T. Gollisch, C. K. Machens, and D. Jaeger. Modeling single-neuron dynamics and computations: a balance of
845 detail and abstraction. *Science*, 314(5796), 2006.

- 846 [94] R. Brette. What is the most realistic single-compartment model of spike initiation? *PLoS Computational Biology*, 11(4):e1004114,
847 2015.
- 848 [95] T. A. Le, A. G. Baydin, R. Zinkov, and F. Wood. Using synthetic data to train neural networks is model-based reasoning. In *2017
849 International Joint Conference on Neural Networks (IJCNN)*, pages 3514–3521. IEEE, 2017.
- 850 [96] J. Chan, V. Perrone, J. Spence, P. Jenkins, S. Mathieson, and Y. Song. A likelihood-free inference framework for population
851 genetic data using exchangeable neural networks. In *Advances in Neural Information Processing Systems*, pages 8594–8605,
852 2018.
- 853 [97] M. G. Blum and O. François. Non-linear regression models for approximate bayesian computation. *Statistics and Computing*, 20
854 (1):63–73, 2010.
- 855 [98] D. B. Rubin et al. Bayesianly justifiable and relevant frequency calculations for the applied statistician. *The Annals of Statistics*,
856 12(4):1151–1172, 1984.
- 857 [99] J. K. Pritchard, M. T. Seielstad, A. Perez-Lezaun, and M. W. Feldman. Population growth of human y chromosomes: a study of y
858 chromosome microsatellites. *Molecular Biology and Evolution*, 16(12):1791–1798, 1999.
- 859 [100] J. Liepe, P. Kirk, S. Filippi, T. Toni, C. P. Barnes, and M. P. Stumpf. A framework for parameter estimation and model selection
860 from experimental data in systems biology using approximate bayesian computation. *Nature Protocols*, 9(2):439, 2014.
- 861 [101] M. U. Gutmann and J. Corander. Bayesian optimization for likelihood-free inference of simulator-based statistical models. *The
862 Journal of Machine Learning Research*, 17(1):4256–4302, 2016.
- 863 [102] O. J. Britton, A. Bueno-Orovio, K. Van Ammel, H. R. Lu, R. Towart, D. J. Gallacher, and B. Rodriguez. Experimentally calibrated
864 population of models predicts and explains intersubject variability in cardiac cellular electrophysiology. *Proceedings of the
865 National Academy of Sciences*, 110(23):E2098–E2105, 2013.
- 866 [103] B. A. Lawson, C. C. Drovandi, N. Cusimano, P. Burrage, B. Rodriguez, and K. Burrage. Unlocking data sets by calibrating
867 populations of models to data density: A study in atrial electrophysiology. *Science Advances*, 4(1):e1701676, 2018.
- 868 [104] S. N. Wood. Statistical inference for noisy nonlinear ecological dynamic systems. *Nature*, 466(7310), 2010.
- 869 [105] R. P. Costa, P. J. Sjostrom, and M. C. Van Rossum. Probabilistic inference of short-term synaptic plasticity in neocortical
870 microcircuits. *Frontiers in Computational Neuroscience*, 7:75, 2013.
- 871 [106] R. Wilkinson. Accelerating abc methods using gaussian processes. In *AISTATS*, 2014.
- 872 [107] E. Meeds and M. Welling. Gps-abc: Gaussian process surrogate approximate bayesian computation. In *Conference on
873 Uncertainty in Artificial Intelligence*, 2014.
- 874 [108] G. Papamakarios, D. Sterratt, and I. Murray. Sequential neural likelihood: Fast likelihood-free inference with autoregressive
875 flows. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 837–848, 2019.
- 876 [109] J.-M. Lueckmann, G. Bassetto, T. Karaletsos, and J. H. Macke. Likelihood-free inference with emulator networks. In F. Ruiz,
877 C. Zhang, D. Liang, and T. Bui, editors, *Proceedings of The 1st Symposium on Advances in Approximate Bayesian Inference*, volume 96
878 of *Proceedings of Machine Learning Research*, pages 32–53, 2019.
- 879 [110] C. Durkan, G. Papamakarios, and I. Murray. Sequential neural methods for likelihood-free inference. *NeurIPS Bayesian Deep
880 Learning Workshop*, 2018.
- 881 [111] S. Barthelmé and N. Chopin. Expectation propagation for likelihood-free inference. *Journal of the American Statistical Association*,
882 109(505):315–333, 2014.
- 883 [112] C. Schröder, L. Lagnado, B. James, and P. Berens. Approximate bayesian inference for a mechanistic model of vesicle release at
884 a ribbon synapse. *BioRxiv*, page 669218, 2019.
- 885 [113] T. A. Le, A. G. Baydin, and F. Wood. Inference compilation and universal probabilistic programming. In *Artificial Intelligence and
886 Statistics*, pages 1338–1348, 2017.
- 887 [114] M. L. Casado, A. G. Baydin, D. M. Rubio, T. A. Le, F. Wood, L. Heinrich, G. Louppte, K. Cranmer, K. Ng, W. Bhimji, et al.
888 Improvements to inference compilation for probabilistic programming in large-scale scientific simulators. *NeurIPS Workshop on
889 Deep Learning for Physical Sciences*, 2017.
- 890 [115] J. Hermans, V. Begy, and G. Louppte. Likelihood-free mcmc with approximate likelihood ratios. *arXiv preprint arXiv:1903.04057*,
891 2019.

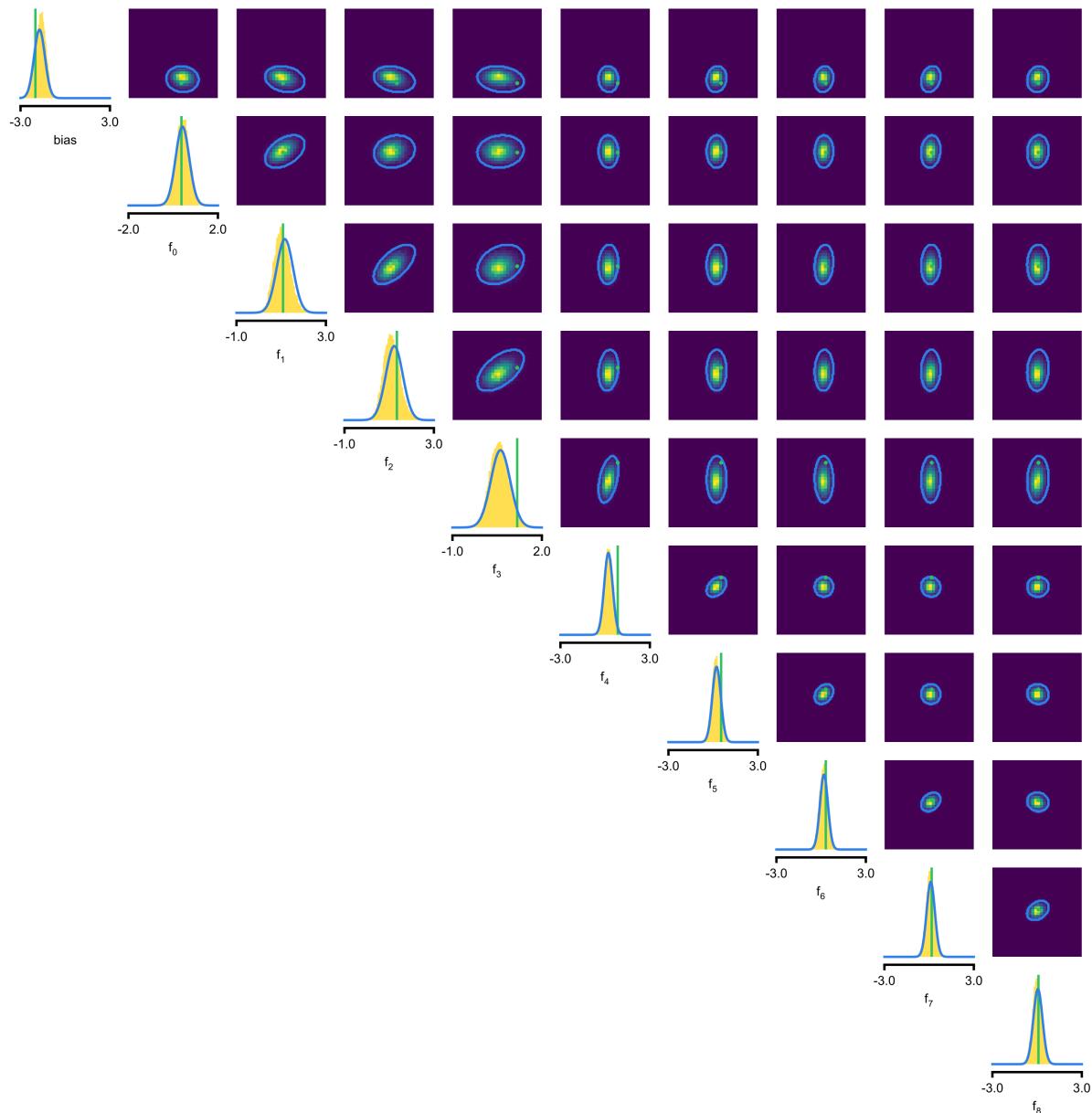
- 892 [116] Z. Chen et al. Bayesian filtering: From kalman filters to particle filters, and beyond. *Statistics*, 182(1):1–69, 2003.
- 893 [117] Q. J. M. Huys and L. Paninski. Smoothing of, and parameter estimation from, noisy biophysical recordings. *PLoS Computational Biology*, 5(5), 2009.
- 894
- 895 [118] L. Hertäg, J. Hass, T. Golovko, and D. Durstewitz. An approximation to the adaptive exponential integrate-and-fire neuron model allows fast and predictive fitting to physiological data. *Frontiers in Computational Neuroscience*, 6:62, 2012.
- 896
- 897 [119] C. Pozzorini, S. Mensi, O. Hagens, R. Naud, C. Koch, and W. Gerstner. Automated high-throughput characterization of single neurons by means of simplified spiking models. *PLoS Computational Biology*, 11(6):e1004275, 2015.
- 898
- 899 [120] J. Ladenbauer, S. McKenzie, D. F. English, O. Hagens, and S. Ostojic. Inferring and validating mechanistic models of neural microcircuits based on spike-train data. *bioRxiv*, page 261016, 2018.
- 900
- 901 [121] A. René, A. Longtin, and J. H. Macke. Inference of a mesoscopic population model from population spike trains. *arXiv preprint arXiv:1910.01618*, 2019.
- 902
- 903 [122] J. Oesterle, C. Behrens, C. Schroeder, T. Herrmann, T. Euler, K. Franke, R. G. Smith, G. Zeck, and P. Berens. Bayesian inference for biophysical neuron models enables stimulus optimization for retinal neuroprosthetics. *bioRxiv*, 2020.
- 904
- 905 [123] S. R. Bittner, A. Palmigiano, A. T. Piet, C. A. Duan, C. D. Brody, K. D. Miller, and J. P. Cunningham. Interrogating theoretical models of neural computation with deep inference. *bioRxiv*, 2019. doi:10.1101/837567.
- 906
- 907 [124] G. Loaiza-Ganem, Y. Gao, and J. P. Cunningham. Maximum entropy flow networks. In *5th International Conference on Learning Representations, ICLR*, 2017.
- 908
- 909 [125] D. Sussillo and L. F. Abbott. Generating coherent patterns of activity from chaotic neural networks. *Neuron*, 63(4):544–557, 2009.
- 910
- 911 [126] V. Mante, D. Sussillo, K. V. Shenoy, and W. T. Newsome. Context-dependent computation by recurrent dynamics in prefrontal cortex. *Nature*, 503(7474):78, 2013.
- 912
- 913 [127] D. Sussillo and O. Barak. Opening the black box: low-dimensional dynamics in high-dimensional recurrent neural networks. *Neural Computation*, 25(3):626–649, 2013.
- 914
- 915 [128] N. Maheswaranathan, A. Williams, M. D. Golub, S. Ganguli, and D. Sussillo. Reverse engineering recurrent networks for sentiment classification reveals line attractor dynamics. *arXiv preprint arXiv:1906.10720*, 2019.
- 916
- 917 [129] D. B. Rubin, S. D. Van Hooser, and K. D. Miller. The stabilized supralinear network: a unifying circuit motif underlying multi-input integration in sensory cortex. *Neuron*, 85(2):402–417, 2015.
- 918
- 919 [130] R. Ratcliff and G. McKoon. The diffusion decision model: theory and data for two-choice decision tasks. *Neural Computation*, 20(4):873–922, 2008.
- 920
- 921 [131] M. G. B. Blum, M. A. Nunes, D. Prangle, S. A. Sisson, et al. A comparative review of dimension reduction methods in approximate bayesian computation. *Statistical Science*, 28(2), 2013.
- 922
- 923 [132] B. Jiang, T.-y. Wu, C. Zheng, and W. H. Wong. Learning summary statistic for approximate bayesian computation via deep neural network. *Statistica Sinica*, pages 1595–1618, 2017.
- 924
- 925 [133] R. Izbicki, A. B. Lee, and T. Pospisil. Abc-cde: Toward approximate bayesian computation with complex high-dimensional data and limited simulations. *Journal of Computational and Graphical Statistics*, pages 1–20, 2019.
- 926
- 927 [134] S. R. Cook, A. Gelman, and D. B. Rubin. Validation of software for bayesian models using posterior quantiles. *Journal of Computational and Graphical Statistics*, 15(3):675–692, 2006.
- 928
- 929 [135] S. Talts, M. Betancourt, D. Simpson, A. Vehtari, and A. Gelman. Validating bayesian inference algorithms with simulation-based calibration. *arXiv preprint arXiv:1804.06788*, 2018.
- 930
- 931 [136] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2014.
- 932
- 933 [137] G. De Nicolao, G. Sparacino, and C. Cobelli. Nonparametric input estimation in physiological systems: problems, methods, and case studies. *Automatica*, 33(5), 1997.
- 934
- 935 [138] L. Abbott and E. Marder. Modeling small networks, 1998.
- 936 [139] J. B. Rosen. The gradient projection method for nonlinear programming. part i. linear constraints. *Journal of the Society for Industrial and Applied Mathematics*, 8(1):181–217, 1960.
- 937

938 **Supplementary material**

939 **Supplementary figures**

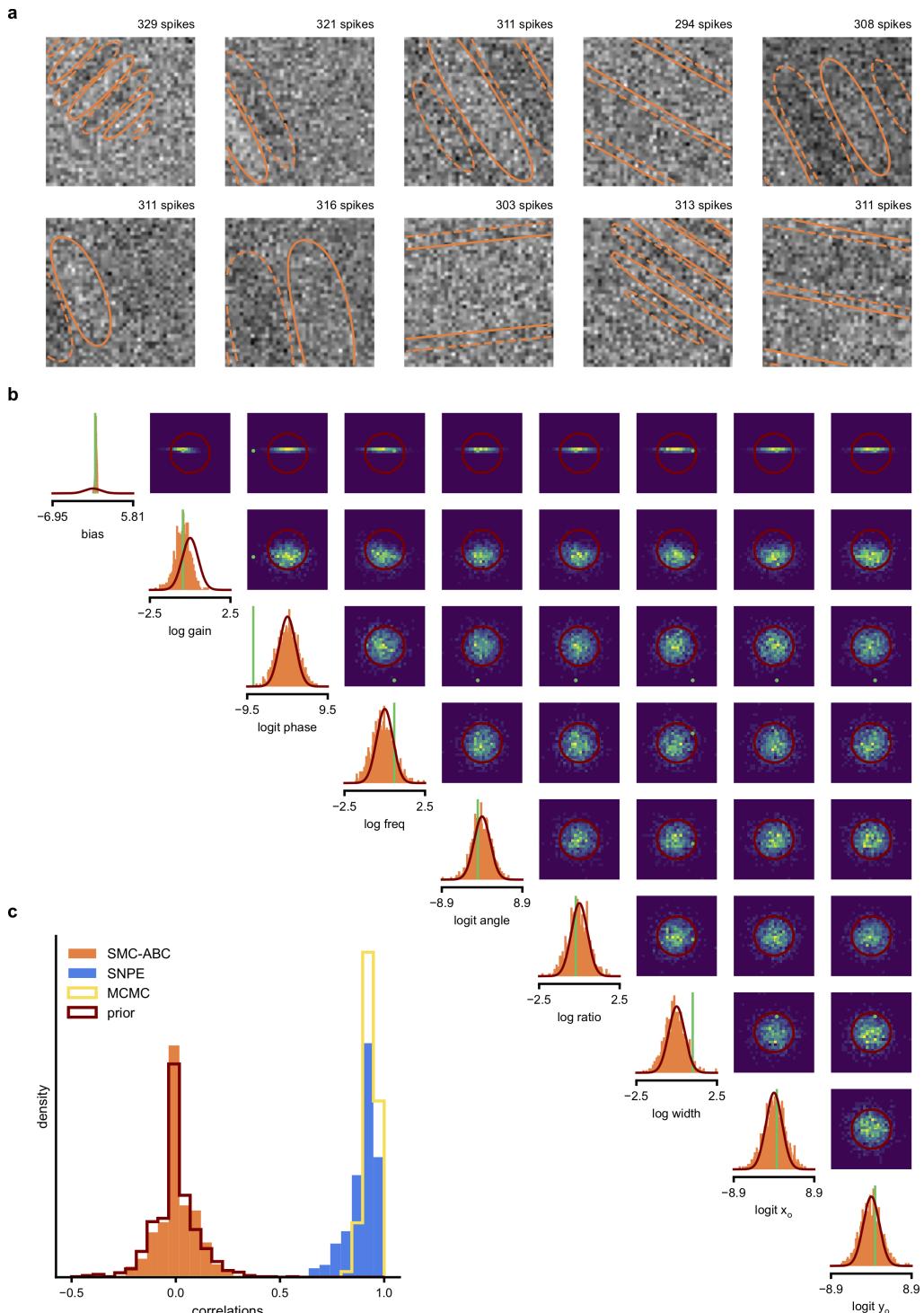


941 **Supplementary Figure 1. Comparison between SNPE-estimated posterior and reference posterior (obtained via**
942 **MCMC) on LN model.** (a) Posterior mean \pm one standard deviation of temporal filter (receptive field) from SNPE posterior
943 (**SNPE**, blue) and reference posterior (MCMC, yellow). (b) Full covariance matrices from SNPE posterior (left) and reference
945 (MCMC, right).

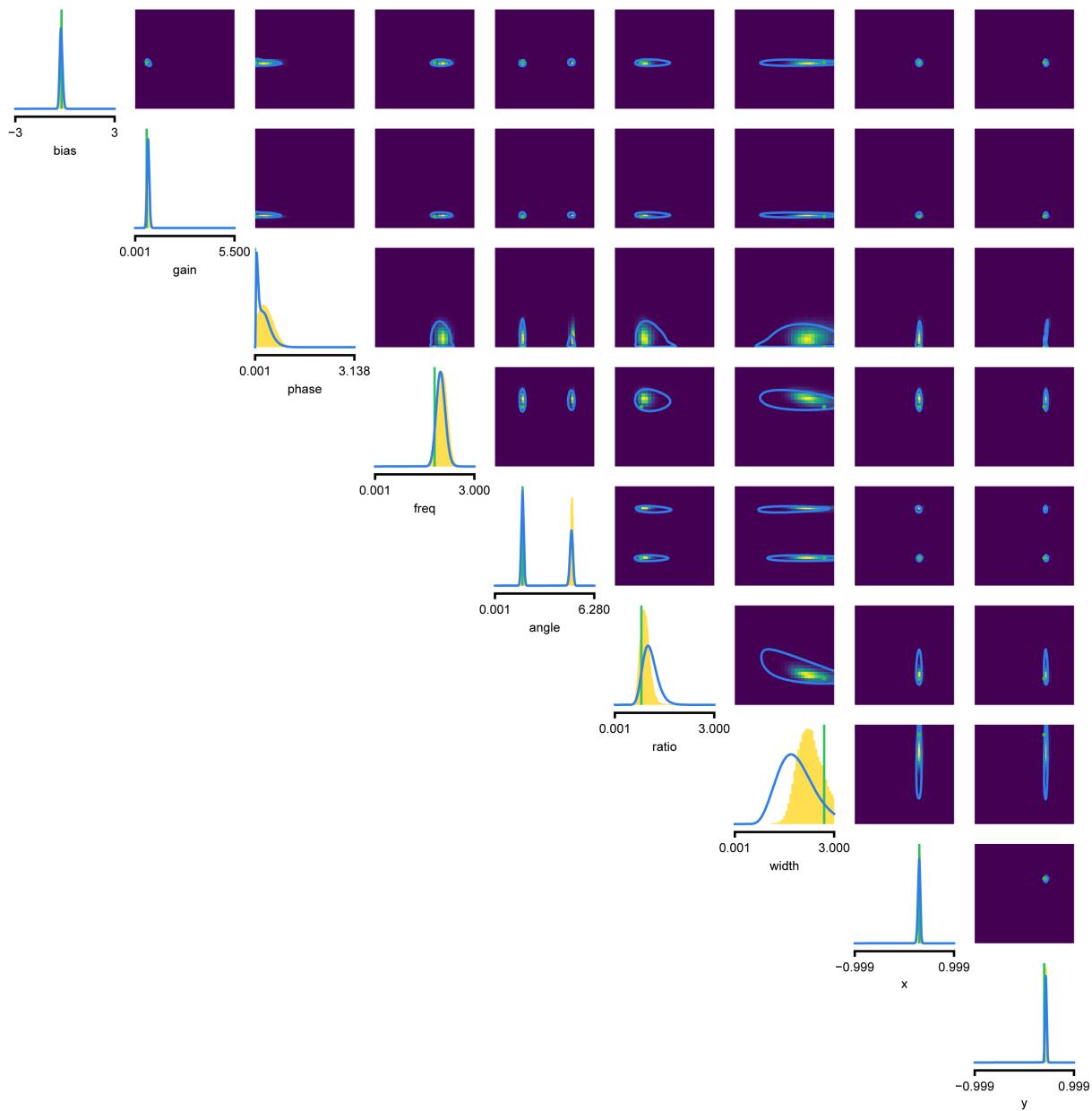


947
948

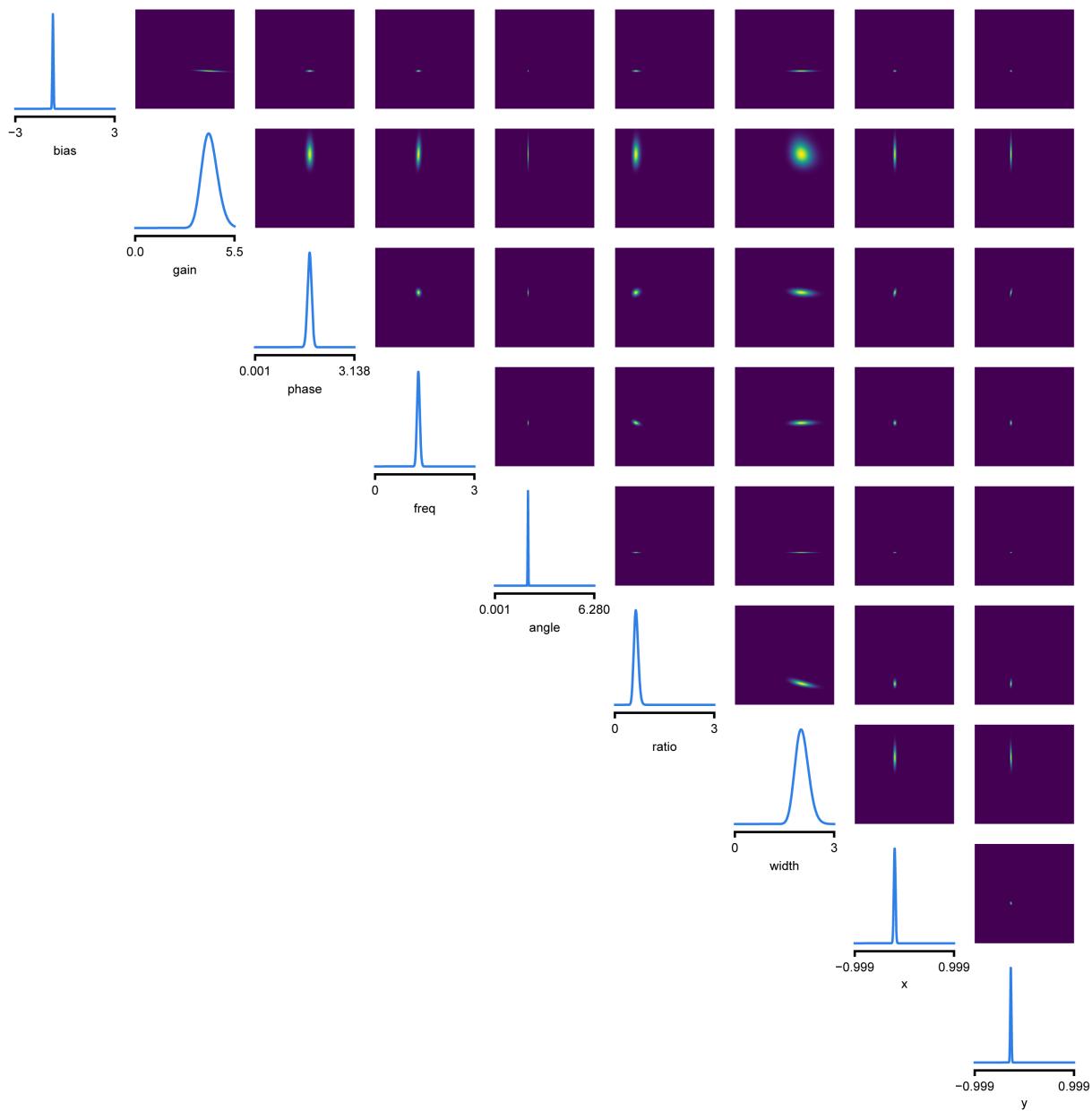
Supplementary Figure 2. Full posterior for LN model. In green, ground-truth parameters. Marginals (blue lines) and 2D marginals for SNPE (contour lines correspond to 95% of the mass) and MCMC (yellow histograms).



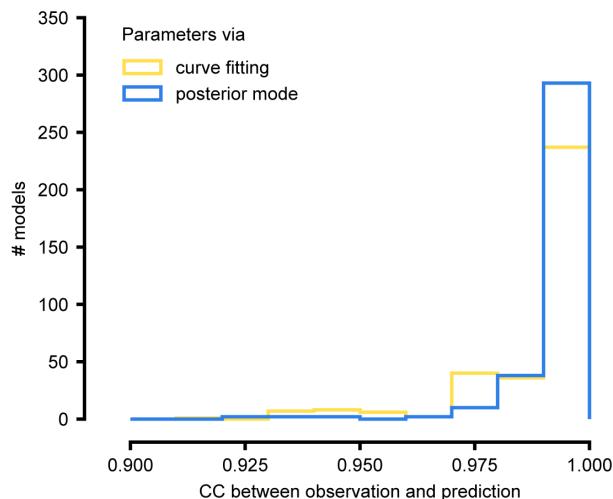
951 **Supplementary Figure 3. SMC-ABC posterior estimate for Gabor GLM receptive field model.** (a) Spike-triggered
952 averages (STAs) and spike counts with closest distance $d(x_o, x_i)$ to the observed data x_o out of 10000 simulations with θ_i
953 sampled from the prior. Spike counts are comparable to the observed data (x_o : 299 spikes), but receptive fields (contours)
954 are not well constrained. (b) Results for SMC-ABC with 10^6 simulations total. Histograms of 1000 particles (orange) returned
955 in the final iteration of SMC-ABC, compared to prior (red contour lines) and ground-truth parameters (green). Distributions
956 over (log-/logit-)transformed parameters, axis limits scaled to mean ± 3 standard deviations of the prior. (c) Correlations
957 between ground-truth receptive field and receptive fields sampled from SMC-ABC posterior (orange), SNPE posterior (blue),
958 reference MCMC posterior (yellow) and prior (red). The SNPE-estimated receptive fields are almost as good as those of the
959 reference posterior, the SMC-ABC estimated ones no better than the prior.



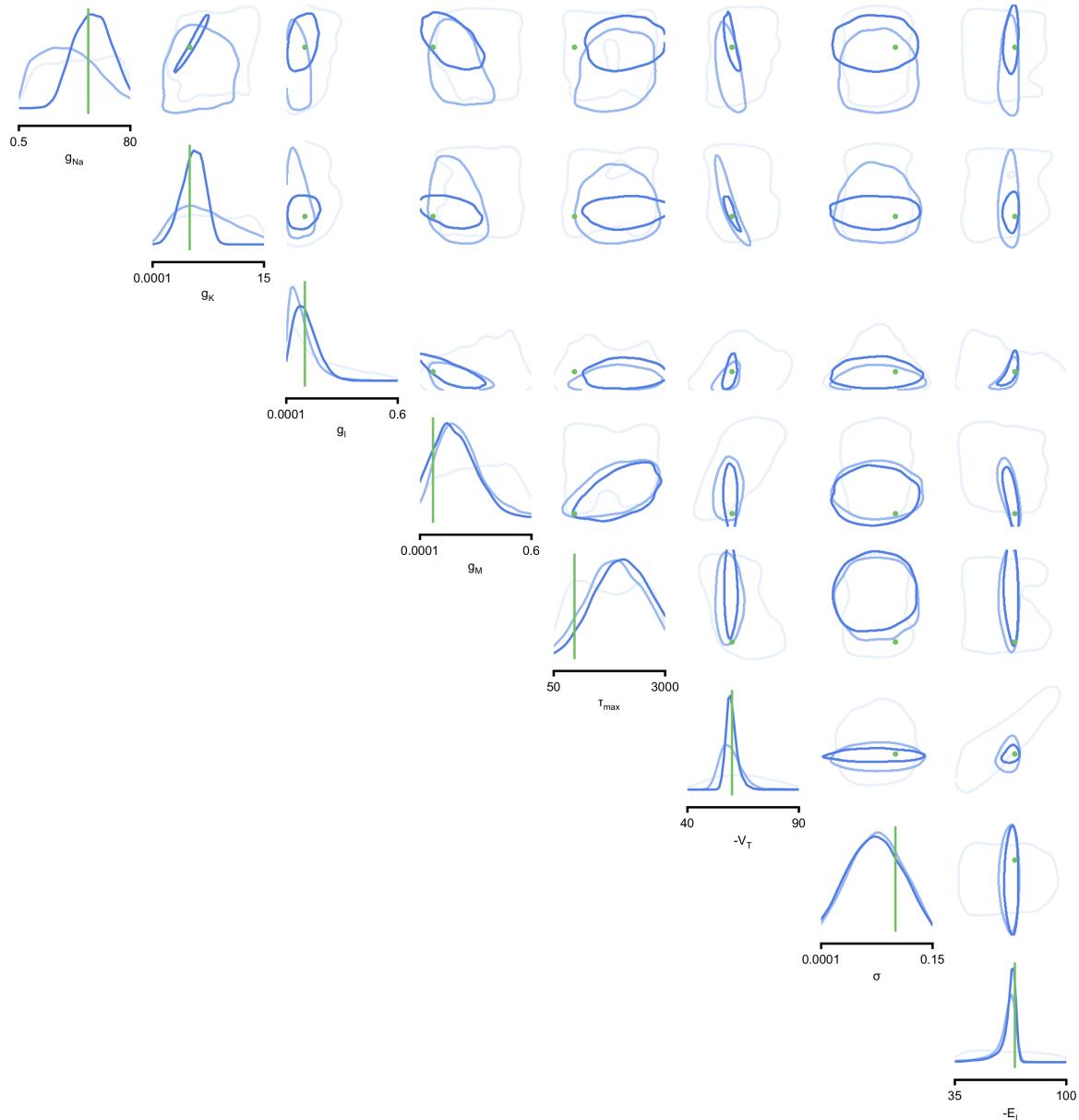
Supplementary Figure 4. Full posterior for Gabor GLM receptive field model. SNPE posterior estimate (blue lines) compared to reference posterior (MCMC, histograms). Ground-truth parameters used to simulate the data in green. We depict the distributions over the original receptive field parameters, whereas we estimate the posterior as a Gaussian mixture over transformed parameters, see Methods for details. We find that a (back-transformed) Gaussian mixture with four components approximates the posterior well in this case.



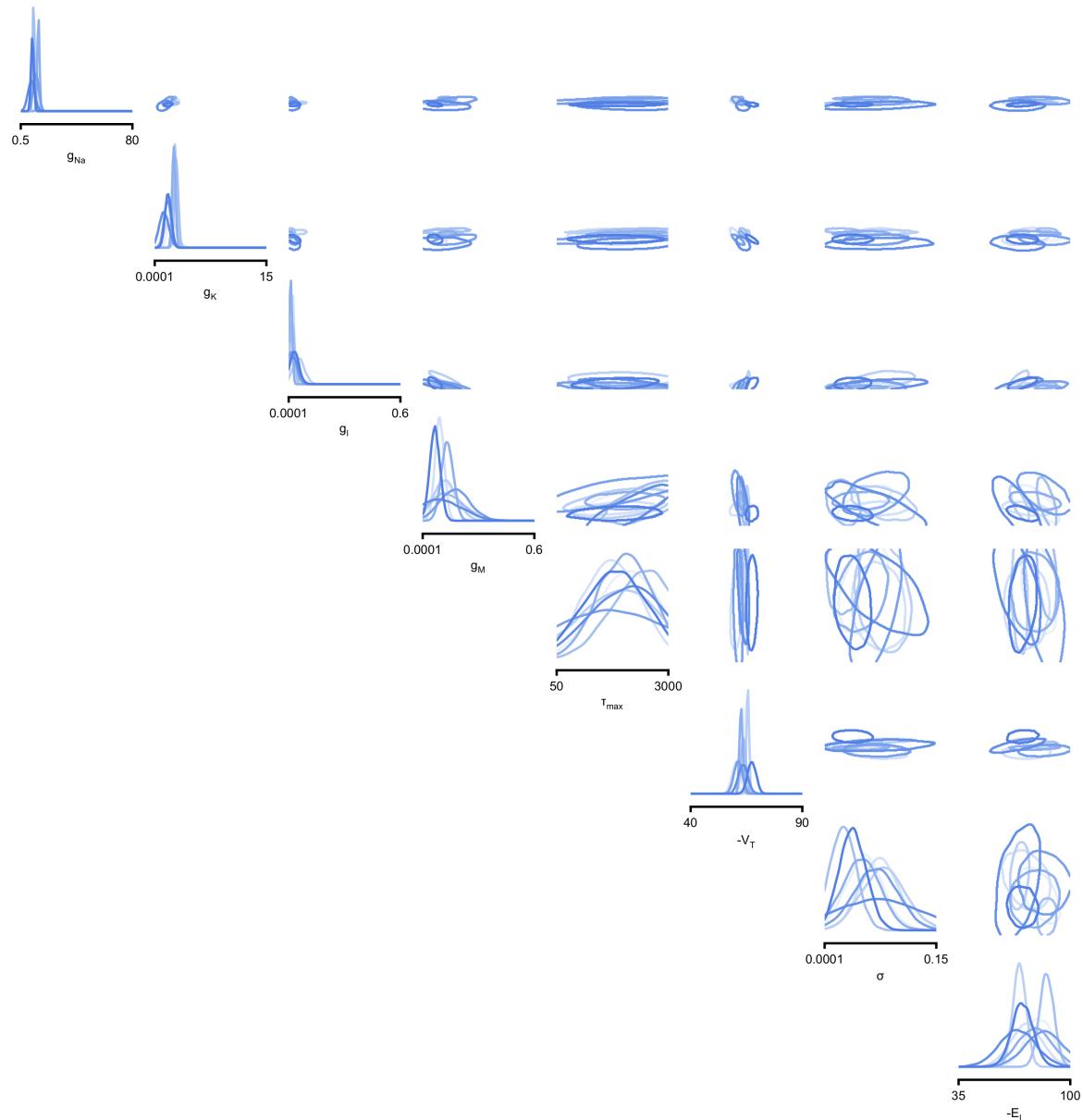
969
970 **Supplementary Figure 5. Full posterior for Gabor LN receptive field model on V1 recordings.** We depict the
971 distributions over the receptive field parameters, derived from the Gaussian mixture over transformed-parameters (see
972 Methods for details).



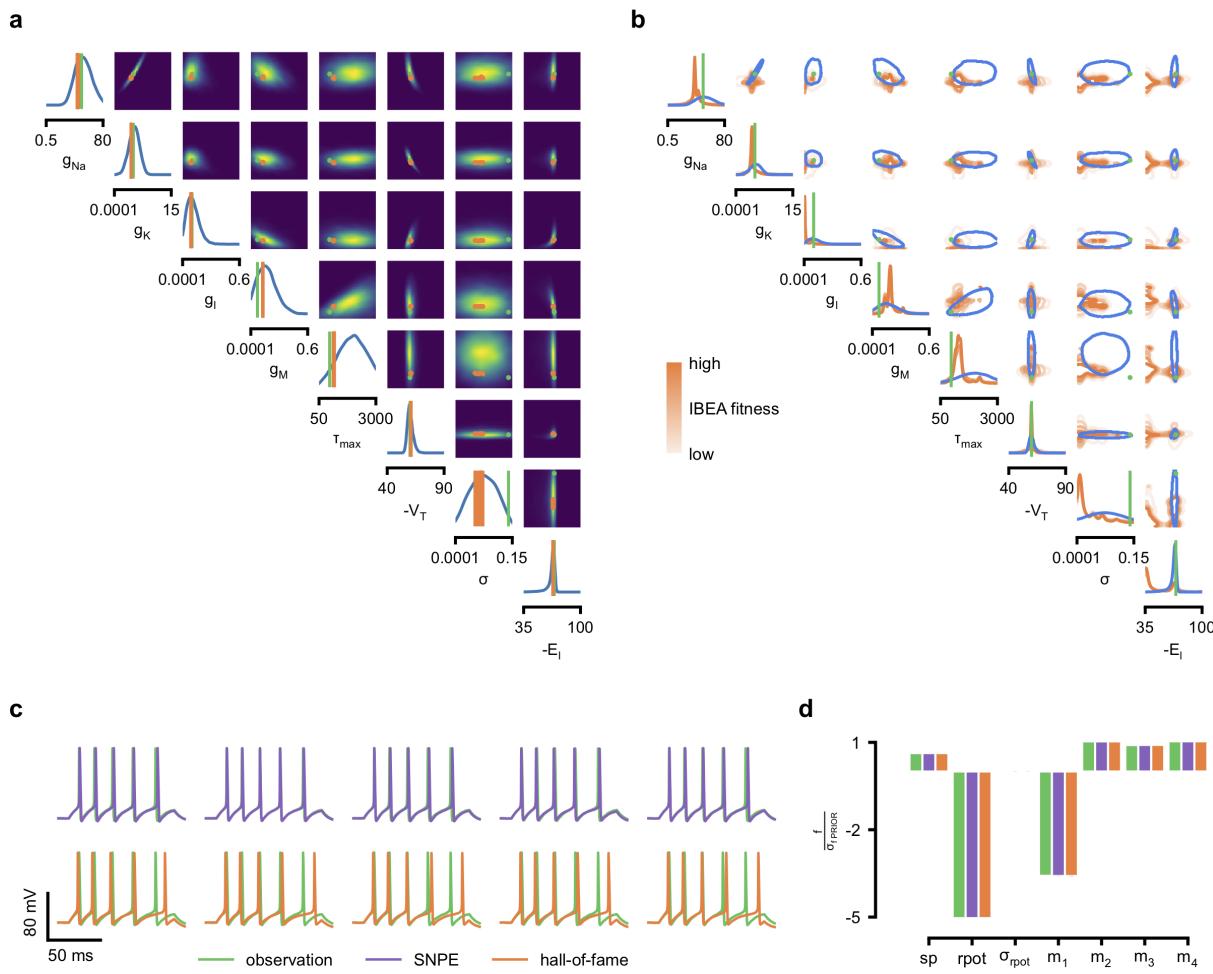
974 **Supplementary Figure 6. Summary results on ICG channel models, and comparison with direct fits.** We generate
975 predictions either with the posterior mode (blue) or with parameters obtained by directly fitting steady-state activation and
976 time-constant curves (yellow). We calculate the correlation coefficient (CC) between observation and prediction. The
977 distribution of CCs is similar for both approaches.
978



980 **Supplementary Figure 7. Full posteriors for Hodgkin-Huxley model for 1, 4 and 7 features.** Images show the pairwise
981 marginals for 7 features. Each contour line corresponds to 68% density mass for a different inferred posterior. Light blue
982 corresponds to 1 feature and dark blue to 7 features.

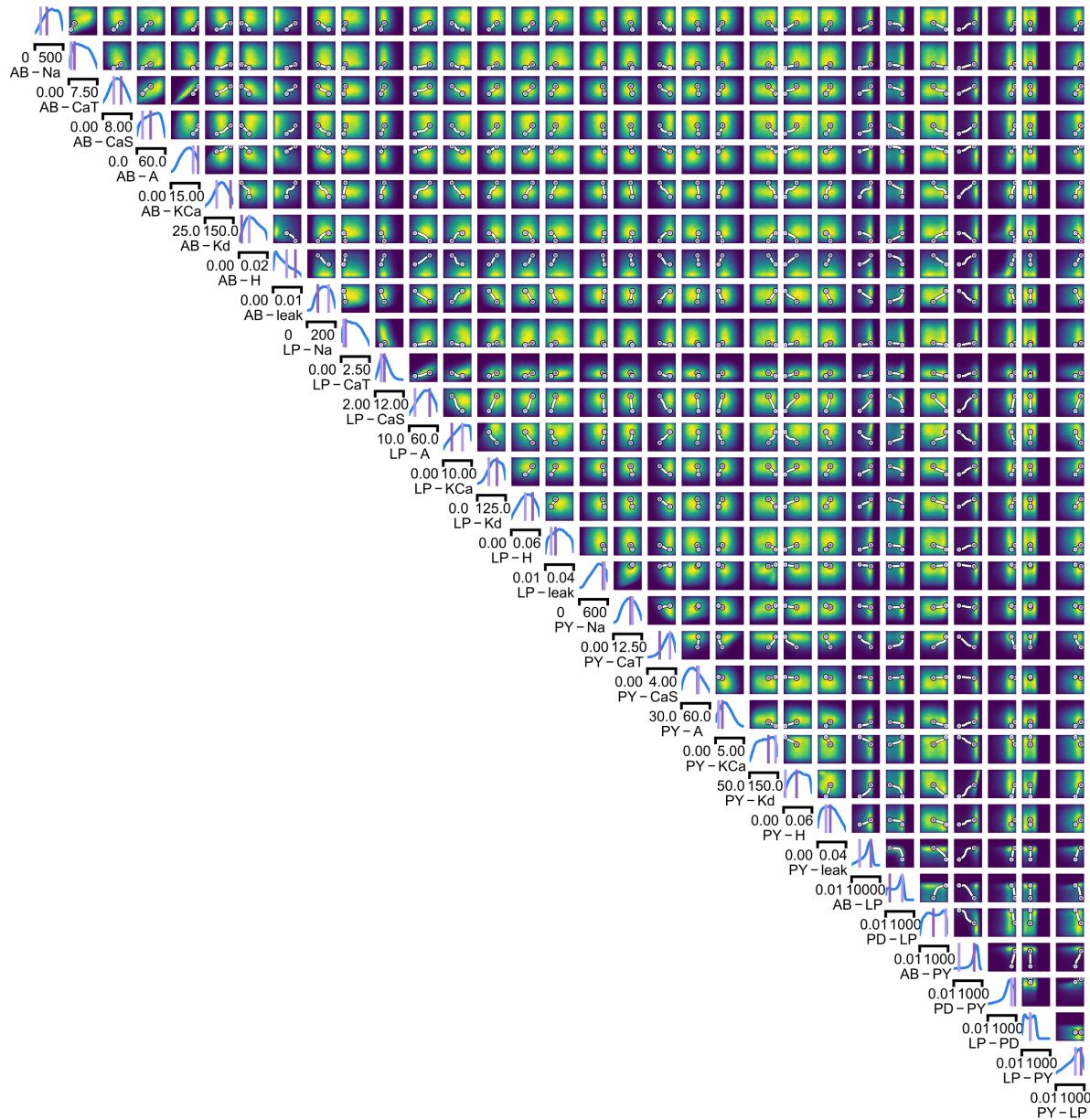


985 **Supplementary Figure 8. Full posteriors for Hodgkin-Huxley model on 8 different recordings from Allen Cell Type
986 Database.** Images show the pairwise marginals for 7 features. Each contour line corresponds to 68% density mass for a
988 different inferred posterior.



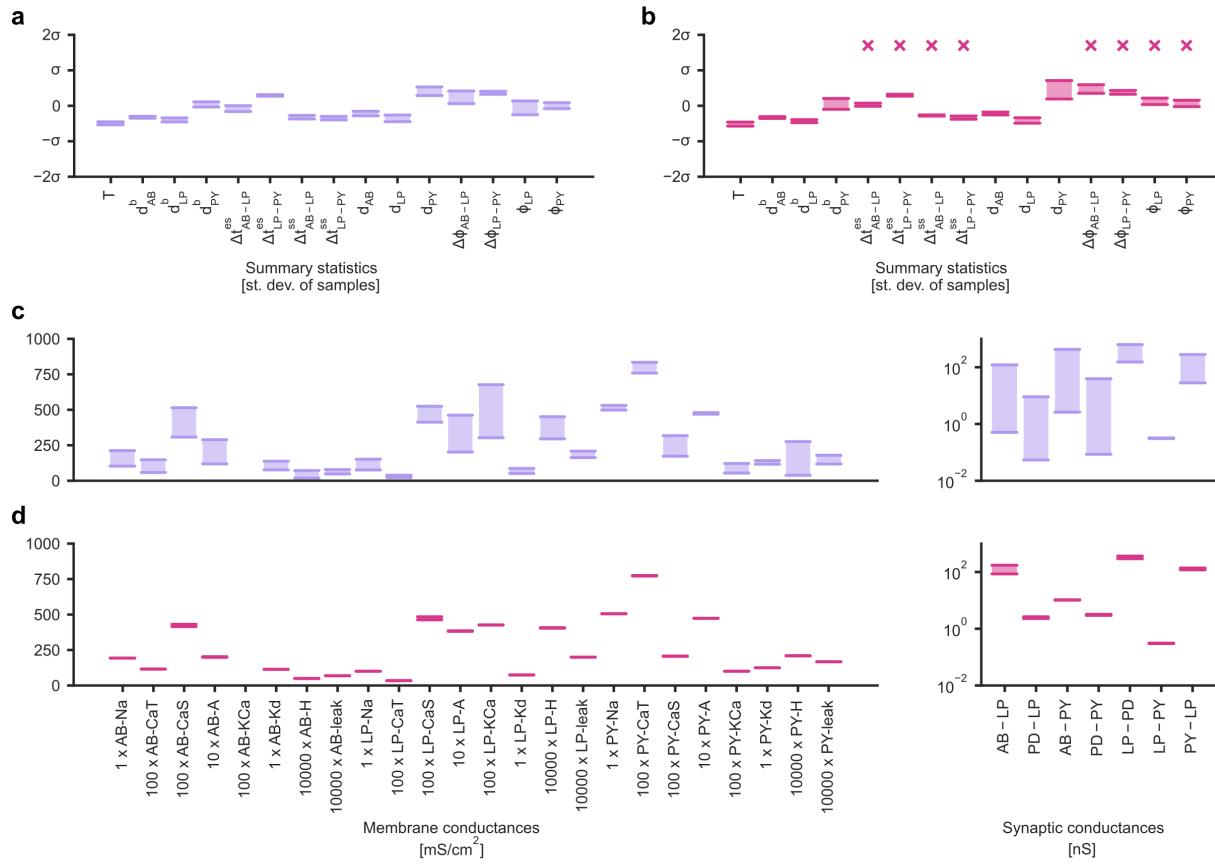
990
991
992
993
994
995
996
997
998
1000

Supplementary Figure 9. Comparison between SNPE posterior and IBEA samples for Hodgkin-Huxley model with 8 parameters and 7 features. (a) Full SNPE posterior distribution. Ground truth parameters in green and IBEA 10 parameters with highest fitness ('hall-of-fame') in orange. (b) Blue contour line corresponds to 68% density mass for SNPE posterior. Light orange corresponds to IBEA sampled parameters with lowest IBEA fitness and dark orange to IBEA sampled parameters with highest IBEA fitness. This plot shows that, in general, SNPE and IBEA can return very different answers- this is not surprising, as both algorithms have different objectives, but this highlights that genetic algorithms do not in general perform statistical inference. (c) Traces for samples with high probability under SNPE posterior (purple), and for samples with high fitness under IBEA objective (hall-of-fame; orange traces). (d) Features for the desired output (observation), the mode of the inferred posterior (purple) and the best sample under IBEA objective (orange). Each voltage feature is normalized by $\sigma_{f \text{ PRIORITY}}$, the standard deviation of the respective feature of simulations sampled from the prior.

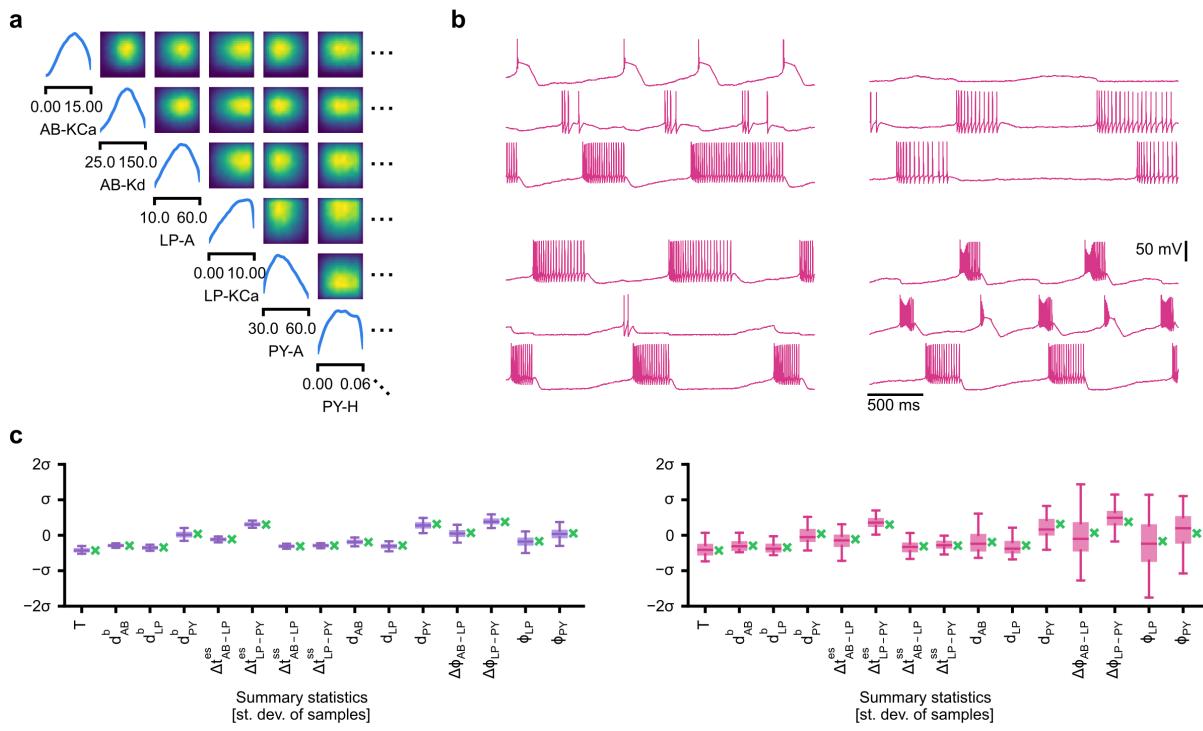


1002
1003
1004
1005

Supplementary Figure 10. Full posterior for the stomatogastric ganglion over 24 membrane and 7 synaptic conductances. The first 24 dimensions depict membrane conductances (top left), the last 7 depict synaptic conductances (bottom right). All synaptic conductances are logarithmically spaced. Between two samples from the posterior with high posterior probability (purple dots), there is a path of high posterior probability (white).



Supplementary Figure 11. Identifying directions of sloppiness and stiffness in the pyloric network of the crustacean stomatogastric ganglion. (a) Minimal and maximal values of all summary statistics along the path lying in regions of high posterior probability, sampled at 20 evenly spaced points. Summary statistics change only little. The summary statistics are scaled with the standard deviation of the 170,000 bursting samples in the created dataset. (b) Summary statistics sampled at 20 evenly spaced points along the orthogonal path. The summary statistics show stronger changes than in panel a and, in particular, often could not be defined because neurons bursted irregularly, as indicated by an 'x' above barplots. (c) Minimal and maximal values of the circuit parameters along the path lying in regions of high posterior probability. Both membrane conductances (left) and synaptic conductances (right) vary over large ranges. Axes as in panel (d). (d) Circuit parameters along the orthogonal path. The difference between the minimal and maximal value is much smaller than in panel (c).



Supplementary Figure 12. Evaluating circuit configurations in which parameters have been sampled independently

(a) Factorized posterior, i.e. posterior obtained by sampling each parameter independently from the associated marginals. Many of the pairwise marginals look similar to the full posterior shown in Supplementary Fig. 10, as the posterior correlations are low. (b) Samples from the factorized posterior- only a minority of these samples produce pyloric activity, highlighting the significance of the posterior correlations between parameters. (c) Left: summary features for 500 samples from the posterior. Boxplot for samples where all summary features are well-defined (80 % of all samples). Right: summary features for 500 samples from the factorized posterior. Only 23 % of these samples have well-defined summary features. The summary features from the factorized posterior have higher variation than the posterior ones. Summary features are normalized using the mean and standard deviation of all samples in our training dataset obtained from prior samples. The boxplots indicate the maximum, 75% quantile, median, 25% quantile, and minimum. The green 'x' indicates the value of the experimental data (the observation, shown in figure 5B).