
BAYESFLOW: LEARNING COMPLEX STOCHASTIC MODELS WITH INVERTIBLE NEURAL NETWORKS

A PREPRINT

Stefan T. Radev

Institute of Psychology
Heidelberg University
Hauptstr. 47-51, 69117 Heidelberg
stefan.radev93@gmail.com

Ulf K. Mertens

Institute of Psychology
Heidelberg University
Hauptstr. 47-51, 69117 Heidelberg
mertens.ulf@gmail.com

Andreas Voss

Institute of Psychology
Heidelberg University
Hauptstr. 47-51, 69117 Heidelberg
andreas.voss@psychologie.uni-heidelberg.de

Lynton Ardizzone

Visual Learning Lab, IWR
Heidelberg University
Im Neuenheimer Feld 205, 69120 Heidelberg
lynton.ardizzone@iwr.uni-heidelberg.de

Ullrich Köthe

Visual Learning Lab, IWR
Heidelberg University
Im Neuenheimer Feld 205, 69120 Heidelberg
ullrich.koethe@iwr.uni-heidelberg.de

December 3, 2020

ABSTRACT

Estimating the parameters of mathematical models is a common problem in almost all branches of science. However, this problem can prove notably difficult when processes and model descriptions become increasingly complex and an explicit likelihood function is not available. With this work, we propose a novel method for globally amortized Bayesian inference based on invertible neural networks which we call BayesFlow. The method uses simulation to learn a global estimator for the probabilistic mapping from observed data to underlying model parameters. A neural network pre-trained in this way can then, without additional training or optimization, infer full posteriors on arbitrary many real datasets involving the same model family. In addition, our method incorporates a summary network trained to embed the observed data into maximally informative summary statistics. Learning summary statistics from data makes the method applicable to modeling scenarios where standard inference techniques with hand-crafted summary statistics fail. We demonstrate the utility of BayesFlow on challenging intractable models from population dynamics, epidemiology, cognitive science and ecology. We argue that BayesFlow provides a general framework for building amortized Bayesian parameter estimation machines for any forward model from which data can be simulated.

1 Introduction

The goal of Bayesian analysis is to infer the underlying characteristics of some natural process of interest given observable manifestations x . In a Bayesian setting, we assume that we already possess sufficient understanding of the forward problem, that is, a suitable model of the mechanism that generates observations from a given configuration of the hidden parameters θ . This forward model can be provided in two forms: In likelihood-based approaches, the likelihood function $p(x | \theta)$ is *explicitly known* and can be *evaluated* analytically or numerically for any pair (x, θ) . In contrast, likelihood-free approaches only require the ability to *sample* from the likelihood. The latter approaches are

typically realized by simulation programs, which generate synthetic observations by means of a deterministic function g of parameters θ and independent noise (i.e., random numbers) ξ :

$$\mathbf{x}_i \sim p(\mathbf{x} | \theta) \iff \mathbf{x}_i = g(\theta, \xi_i) \text{ with } \xi_i \sim p(\xi) \quad (1)$$

In such cases, the likelihood $p(\mathbf{x} | \theta)$ is only defined *implicitly* via the action of the simulation program g , but calculation of its actual numerical value for a simulated observation \mathbf{x}_i is impossible. This, in turn, prohibits standard statistical inference.

Likelihood-free problems arise, for example, when $p(\mathbf{x} | \theta)$ is not available in closed-form, or when the forward model is defined by a stochastic differential equation, a Monte-Carlo simulation, or a complicated algorithm [27, 49, 47, 51]. In this paper, we propose a new Bayesian solution to the likelihood-free setting in terms of *invertible neural networks*.

Bayesian modeling leverages the available knowledge about the forward model to get the best possible estimate of the posterior distribution of the inverse model:

$$p(\theta | \mathbf{x}_{1:N}) = \frac{p(\mathbf{x}_{1:N} | \theta) p(\theta)}{\int p(\mathbf{x}_{1:N} | \theta) p(\theta) d\theta} \quad (2)$$

In Bayesian inference, the posterior encodes all information about θ obtainable from a set of observations $\mathbf{x}_{1:N} = \{\mathbf{x}_i\}_{i=1}^N$. The observations are assumed to arise from N runs of the forward model with fixed, but unknown, true parameters θ^* . Bayesian inverse modeling is challenging for three reasons:

1. The right-hand side of Bayes' formula above is always intractable in the likelihood-free case and must be approximated.
2. The forward model is usually non-deterministic, so that there is intrinsic uncertainty about the true value of θ .
3. The forward model is typically not information-preserving, so that there is ambiguity among possible values of θ .

The standard solution to these problems is offered by *approximate Bayesian computation* (ABC) methods [45, 10, 39, 47]. ABC methods approximate the posterior by repeatedly sampling parameters from a proposal (prior) distribution $\theta^{(l)} \sim p(\theta)$ and then simulating multiple datasets by running the forward model $\mathbf{x}_i \sim p(\mathbf{x} | \theta^{(l)})$ for $i = 1 \dots N$. If the resulting dataset is sufficiently similar to the actually observed dataset $\mathbf{x}_{1:N}^o$, the corresponding $\theta^{(l)}$ is retained as a sample from the desired posterior, otherwise rejected. Stricter similarity criteria lead to more accurate approximations of the desired posterior at the price of higher and oftentimes prohibitive rejection rates.

More efficient methods for approximate inference, such as sequential Monte Carlo (ABC-SMC), Markov-Chain Monte Carlo variants [44], or the recent neural density estimation methods [16, 38, 30], optimize sampling from a proposal distribution in order to balance the speed-accuracy trade-off of vanilla ABC methods. More details can be found in the section **Related Work** and in the excellent review by [9].

All sampling methods described above operate on the level of individual datasets, that is, for each observation sequence $\mathbf{x}_{1:N}$, the entire estimation procedure for the posterior must be run again from scratch. Therefore, we refer to this approach as *case-based inference*. Running estimation for each individual dataset separately stands in contrast to *amortized inference*, where estimation is split into a potentially expensive *upfront* training phase, followed by a much cheaper inference phase. The goal of the upfront training phase is to learn an approximate posterior $\hat{p}(\theta | \mathbf{x}_{1:N})$ that works well for *any* observation sequence $\mathbf{x}_{1:N}$. Evaluating this model for specific observations $\mathbf{x}_{1:N}^o$ is then very fast, so that the training effort amortizes over repeated evaluations (see Figure 1 for a graphical illustration). The break-even between case-based and amortized inference depends on the application and model types, and we will report comparisons in the experimental section. Our main aim in this paper, however, is the introduction of a general approach to amortized Bayesian inference and the demonstration of its excellent accuracy in posterior estimation for a variety of popular forward models.

To make amortized inference feasible in practice, it must work well for arbitrary dataset sizes N . Depending on data acquisition circumstances, the number of available observations for a fixed model parameter setting may vary considerably, ranging from $N = 1$ to several hundreds and beyond. This has not only consequences for the required architecture of our density approximators, but also for their behavior: They must exhibit correct *posterior contraction*. Accordingly, the estimated posterior $\hat{p}(\theta | \mathbf{x}_{1:N})$ should get sharper (i.e., more peaked) as the number N of available observations increases. In the simplest case, the posterior variance should decrease at rate $1/N$, but more complex behavior can occur for difficult (e.g., multi-modal) true posteriors $p(\theta | \mathbf{x}_{1:N})$.

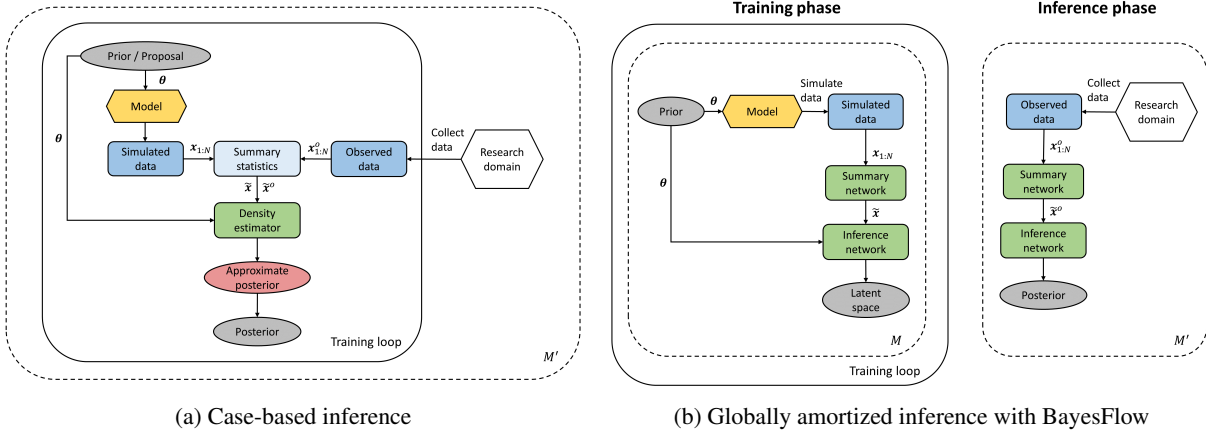


Figure 1: Graphical illustration of the main differences between case-based (neural) density estimation methods and BayesFlow. **(a)** Case-based methods require a separate optimization loop for each observed dataset from a given research domain. When case-based methods incorporate a training phase (e.g., APT), it must be repeated for each new dataset. Summary statistics are manually selected and may thus be sub-optimal; **(b)** BayesFlow incorporates a global upfront training (before any real data are collected) via simulations from the forward model (left panel). Summary and inference network are trained jointly, resulting in higher accuracy than hand-crafted summary statistics. In the inference phase (right panel), BayesFlow works entirely in a feed-forward manner, that is, no training or optimization happens in this phase. The upfront training effort is therefore amortized over arbitrary many observed datasets from a research domain working on the same model family. Note that the solid and dashed plates are swapped between case-based Bayesian inference and the training phase of BayesFlow.

We incorporate these considerations into our method by integrating two separate deep neural networks modules (detailed in the **Methods** section; see also Figure 1), which are trained jointly on simulated data from the forward model: a *summary network* and an *inference network*.

The *summary network* is responsible for reducing a set of observations $x_{1:N}$ of variable size to a fixed-size vector of *learned* summary statistics. In traditional likelihood-free approaches, the method designer is responsible for selecting suitable statistics for each application *a priori* [33, 32, 43, 45]. In contrast, our summary networks learn the most informative statistics directly from data, and we will show experimentally (see **Experiment 3.8**) that these statistics are superior to manually constructed ones. Summary networks differ from standard feed-forward networks because they should be independent of the input size N and respect the inherent functional and probabilistic symmetries of the data. For example, permutation invariant networks are required for *i.i.d.* observations [6], and recurrent networks [15] or convolutional networks [29] for data with temporal or spatial dependencies.

The *inference network* is responsible for learning the true posterior of model parameters given the summary statistics of the observed data. Since it sees the data only through the lens of the summary network, all symmetries captured by the latter are automatically inherited by the posterior. We implement the inference network as an *invertible neural network*. Invertible neural networks are based on the recent theory and applications of *normalizing flows* [3, 25, 18, 13, 26]. Flow-based methods can perform exact inference under perfect convergence and scale favourably from simple low-dimensional problems to high-dimensional distributions with complex dependencies, for instance, the pixels of an image [25]. For each application/model of interest, we train an invertible network jointly with a corresponding summary network using simulated data from the respective known forward model with reasonable priors. After convergence of this forward training, the network’s invertibility ensures that a model for the inverse model is obtained for free, simply by running inference through the model backwards. Thus, our networks can perform fast amortized Bayesian inference on arbitrary many datasets from a given application domain without expensive case-based optimization. We call our method *BayesFlow*, as it combines ideas from Bayesian inference and flow-based deep learning.

BayesFlow draws on major advances in modern deep probabilistic modeling, also referred to as deep generative modeling [6, 25, 2, 24]. A hallmark idea in deep probabilistic modeling is to represent a complicated target distribution as a non-linear bijective transformation of some simpler latent distribution (e.g., Gaussian or uniform), a so called *pushforward*. Density estimation of the target distribution, a very complex problem, is thus reduced to learning a non-linear transformation, a task that is ideally suited for gradient-based neural network training via standard backpropagation. During the *inference phase*, samples from the target distribution are obtained by sampling from the simpler latent distribution and applying the inverse transformation learned during the training phase (see Figure 1b

for a high-level overview). Using this approach, recent applications of deep probabilistic models have achieved unprecedented performance on hitherto intractable high-dimensional problems [6, 25, 18].

In the context of Bayesian inference, the target distribution is the posterior $p(\theta | x_{1:N})$ of model parameters given observed data. We leverage the fact that we can simulate arbitrarily large amounts of training data from the forward model in order to ensure that the summary and invertible networks approximate the true posterior as well as possible. During the inference phase, our model can either numerically evaluate the posterior probability of any candidate parameter θ , or can generate a posterior sample $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(L)}$ of likely parameters for the observed data $x_{1:N}^o$. In the **Methods** section, we show that our networks indeed sample from the correct posterior under perfect convergence. In summary, the contributions of our BayesFlow method are the following:

- Globally amortized approximate Bayesian inference with invertible neural networks;
- Learning maximally informative summary statistics from raw datasets with variable number of observations instead of relying on restrictive hand-crafted summary statistics;
- Theoretical guarantee for sampling from the true posterior distribution with arbitrary priors and posteriors;
- Parallel computations applicable to both forward simulations and neural network optimization;

To illustrate the utility of BayesFlow, we first apply it to two toy models with analytically tractable posteriors. The first is a multivariate Gaussian with a full covariance matrix and a unimodal posterior. The second is a Gaussian mixture model with a multimodal posterior. Then, we present applications to challenging models with intractable likelihoods from population dynamics, cognitive science, epidemiology, and ecology and demonstrate the utility of BayesFlow in terms of speed, accuracy of recovery, and probabilistic calibration. Alongside, we introduce several performance validation tools.

1.1 Related Work

BayesFlow incorporates ideas from previous machine learning and deep learning approaches to likelihood-free inference [31, 41, 33, 43, 22]. The most common approach has been to cast the problem of parameter estimation as a supervised learning task. In this setting, a large dataset of the form $D = \{(h(x_{1:N}^{(m)}), \theta^{(m)})\}_{m=1}^M$ is created by repeatedly sampling from $p(\theta)$ and simulating an artificial datasets $x_{1:N}$ by running the simulator with the sampled parameters. Usually, the dimensionality of the simulated data is reduced by computing summary statistics with a fixed summary function $h(x_{1:N})$. Then, a supervised learning algorithm (e.g., random forest [43], or a neural network [41]) is trained on the summary statistics of the simulated data to output an estimate of the true data generating parameters. Thus, an attempt is made to approximate the intractable inverse model $\theta = g^{-1}(x, \xi)$. A main shortcoming of supervised approaches is that they provide only limited information about the posterior (e.g., point-estimates, quantiles or variance estimates) or impose overly restrictive distributional assumptions on the shape of the posterior (e.g., Gaussian).

Our ideas are also closely related to the concept of *optimal transport maps* and its application in Bayesian inference [12, 40, 8, 5]. A transport map defines a transformation between (probability) measures which can be constructed in a way to *warp* a simple probability distribution into a more complex one. In the context of Bayesian inference, transport maps have been applied to accelerate MCMC sampling [40], to perform sequential inference [12], and to solve inference problems via direct optimization [5]. In fact, BayesFlow can be viewed as a parameterization of invertible transport maps via invertible neural networks. An important distinction is that BayesFlow does not require an explicit likelihood function for approximating the target posteriors and is capable of amortized inference.

Similar ideas for likelihood-free inference are incorporated in the recent automatic posterior transformation (APT) [16], and the sequential neural likelihood (SNL) [38] methods. APT iteratively refines a proposal distribution via masked autoregressive flow (MAF) networks to generate parameter samples which closely match a particular observed dataset. SNL, in turn, trains a masked autoencoder density estimator (MADE) neural network within an MCMC loop to speed-up convergence to the true posterior. Even though these methods also entail a relatively expensive learning phase and a cheap inference phase, posterior inference is amortized only for a single dataset. Thus, the learning phase needs to be run through for each individual dataset (see Figure 1a). In contrast, we propose to learn the posterior globally over the entire range of plausible parameters and datasets by employing a conditional invertible neural network (cINN) estimator (see Figure 1b). Previously, INNs have been successfully employed to model data from astrophysics and medicine [2]. We adapt the model to suit the task of parameter estimation in the context of mathematical modeling and develop a probabilistic architecture for performing fully Bayesian and globally amortized inference with complex mathematical models.

2 Methods

2.1 Notation

In the following, the number of parameters of a mathematical model will be denoted as D , and the number of observations in a dataset as N . We denote data simulated from the mathematical model of interest as $\mathbf{x}_{1:N} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$, where each individual \mathbf{x}_i can represent a scalar or a vector. Observed or test data will be marked with a superscript o (i.e., $\mathbf{x}_{1:N}^o$). The parameters of a mathematical model are represented as a vector $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_D)$, and all trainable parameters of the invertible and summary neural networks as ϕ and ψ , respectively. When a dataset consists of observations over a period of time, the number of observations will be denoted as T .

2.2 Learning the Posterior

Assume that we have an invertible function $f_\phi : \mathbb{R}^D \rightarrow \mathbb{R}^D$, parameterized by a vector of parameters ϕ , for which the inverse $f_\phi^{-1} : \mathbb{R}^D \rightarrow \mathbb{R}^D$ exists. For now, consider the case when raw simulated data $\mathbf{x}_{1:N}$ of size $N = 1$ is entered directly into the invertible network without using a summary network. Our goal is to train an invertible neural network which approximates the true posterior as accurately as possible:

$$p_\phi(\boldsymbol{\theta} | \mathbf{x}) \approx p(\boldsymbol{\theta} | \mathbf{x}) \quad (3)$$

for all possible $\boldsymbol{\theta}$ and \mathbf{x} . We reparameterize the approximate posterior p_ϕ in terms of a conditional invertible neural network (cINN) f_ϕ which implements a normalizing flow between $\boldsymbol{\theta}$ and a Gaussian latent variable \mathbf{z} :

$$\boldsymbol{\theta} \sim p_\phi(\boldsymbol{\theta} | \mathbf{x}) \iff \boldsymbol{\theta} = f_\phi^{-1}(\mathbf{z}; \mathbf{x}) \text{ with } \mathbf{z} \sim \mathcal{N}_D(\mathbf{z} | \mathbf{0}, \mathbb{I}) \quad (4)$$

Accordingly, we need to ensure that the outputs of $f_\phi^{-1}(\mathbf{z}; \mathbf{x})$ follow the target posterior $p(\boldsymbol{\theta} | \mathbf{x})$. Thus, we seek neural network parameters $\hat{\phi}$ which minimize the Kullback-Leibler (KL) divergence between the true and the model-induced posterior for all possible datasets \mathbf{x} . Therefore, our objective becomes:

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} \mathbb{E}_{p(\mathbf{x})} [\mathbb{KL}(p(\boldsymbol{\theta} | \mathbf{x}) || p_\phi(\boldsymbol{\theta} | \mathbf{x}))] \quad (5)$$

$$= \underset{\phi}{\operatorname{argmin}} \mathbb{E}_{p(\mathbf{x})} [\mathbb{E}_{p(\boldsymbol{\theta} | \mathbf{x})} [\log p(\boldsymbol{\theta} | \mathbf{x}) - \log p_\phi(\boldsymbol{\theta} | \mathbf{x})]] \quad (6)$$

$$= \underset{\phi}{\operatorname{argmax}} \mathbb{E}_{p(\mathbf{x})} [\mathbb{E}_{p(\boldsymbol{\theta} | \mathbf{x})} [\log p_\phi(\boldsymbol{\theta} | \mathbf{x})]] \quad (7)$$

$$= \underset{\phi}{\operatorname{argmax}} \iint p(\mathbf{x}, \boldsymbol{\theta}) \log p_\phi(\boldsymbol{\theta} | \mathbf{x}) d\mathbf{x} d\boldsymbol{\theta} \quad (8)$$

Note, that the log posterior density $p(\boldsymbol{\theta} | \mathbf{x})$ can be dropped from the optimization objective in Eq.7, as it does not depend on the neural network parameters ϕ . In other words, we seek neural network parameters $\hat{\phi}$ which maximize the posterior probability of data-generating parameters $\boldsymbol{\theta}$ given observed data \mathbf{x} for all $\boldsymbol{\theta}$ and \mathbf{x} . Since $f_\phi(\boldsymbol{\theta}; \mathbf{x}) = \mathbf{z}$ by design, the change of variable rule of probability yields:

$$p_\phi(\boldsymbol{\theta} | \mathbf{x}) = p(\mathbf{z} = f_\phi(\boldsymbol{\theta}; \mathbf{x})) \left| \det \left(\frac{\partial f_\phi(\boldsymbol{\theta}; \mathbf{x})}{\partial \boldsymbol{\theta}} \right) \right| \quad (9)$$

Thus, we can re-write our objective as:

$$\hat{\phi} = \underset{\phi}{\operatorname{argmax}} \iint p(\mathbf{x}, \boldsymbol{\theta}) \log p_\phi(\boldsymbol{\theta} | \mathbf{x}) d\mathbf{x} d\boldsymbol{\theta} \quad (10)$$

$$= \underset{\phi}{\operatorname{argmax}} \iint p(\mathbf{x}, \boldsymbol{\theta}) (\log p(f_\phi(\boldsymbol{\theta}; \mathbf{x})) + \log |\det \mathbf{J}_{f_\phi}|) d\mathbf{x} d\boldsymbol{\theta} \quad (11)$$

where we have abbreviated $\partial f_\phi(\boldsymbol{\theta}; \mathbf{x}) / \partial \boldsymbol{\theta}$ (the Jacobian of f_ϕ evaluated at $\boldsymbol{\theta}$ and \mathbf{x}) as \mathbf{J}_{f_ϕ} . Due to the architecture of our cINN, the log $|\det \mathbf{J}_{f_\phi}|$ is easy to compute (see next section for details).

Utilizing simulations from the forward model (Eq.1), we can approximate the expectations by minimizing the Monte-Carlo estimate of the negative of Eq.11. Accordingly, for a batch of M simulated datasets and data-generating

parameters $\{(\mathbf{x}^{(m)}, \boldsymbol{\theta}^{(m)})\}_{m=1}^M$ we have:

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} \frac{1}{M} \sum_{m=1}^M -\log p_{\phi}(\boldsymbol{\theta}^{(m)} | \mathbf{x}^{(m)}) \quad (12)$$

$$= \underset{\phi}{\operatorname{argmin}} \frac{1}{M} \sum_{m=1}^M \left(-\log p(f_{\phi}(\boldsymbol{\theta}^{(m)}; \mathbf{x}^{(m)})) - \log |\det \mathbf{J}_{f_{\phi}}^{(m)}| \right) \quad (13)$$

$$= \underset{\phi}{\operatorname{argmin}} \frac{1}{M} \sum_{m=1}^M \left(\frac{\|f_{\phi}(\boldsymbol{\theta}^{(m)}; \mathbf{x}^{(m)})\|_2^2}{2} - \log |\det \mathbf{J}_{f_{\phi}}^{(m)}| \right) \quad (14)$$

We treat Eq.14 as a loss function $\mathcal{L}(\phi)$ which can be minimized with any stochastic gradient descent method. The first term follows from Eq.13 due to the fact that we have prescribed a unit Gaussian distribution to \mathbf{z} . It represents the negative log of $\mathcal{N}_D(\mathbf{z} | \mathbf{0}, \mathbb{I}) \propto \exp(-\frac{1}{2}\|\mathbf{z}\|_2^2)$. The second term controls the rate of volume change induced by the learned non-linear transformation from $\boldsymbol{\theta}$ to \mathbf{z} achieved by f_{ϕ} . Thus, minimizing Eq.14 ensures that \mathbf{z} follows the prescribed unit Gaussian.

The correctness of the learned posterior can be guaranteed in the following way, assuming the network is able to reach the global minimum of the loss (i.e. under perfect convergence).

Proposition 1. *Assume that the cINN architecture and domain of ϕ are chosen such that $\hat{\phi}$ is the global minimum of the objective in Eq.11. Then, the latent output distribution will be statistically independent of the conditioning data, $p_{\hat{\phi}}(\mathbf{z} | \mathbf{x}) \perp p(\mathbf{x})$. As a result, the samples transformed backwards from $p(\mathbf{z})$ will follow the true posterior, that is:*

$$f_{\hat{\phi}}^{-1}(\mathbf{z}; \mathbf{x}) \sim p(\boldsymbol{\theta} | \mathbf{x}) \quad \text{with} \quad \mathbf{z} \sim \mathcal{N}_D(\mathbf{z} | \mathbf{0}, \mathbb{I}) \quad (15)$$

Proof. For short, we denote $p(\mathbf{z}) := \mathcal{N}_D(\mathbf{z} | \mathbf{0}, \mathbb{I})$, and the distribution of network outputs as $p(f_{\hat{\phi}}(\boldsymbol{\theta}; \mathbf{x})) := p_{\hat{\phi}}(\mathbf{z} | \mathbf{x})$. Due to $\mathbb{KL}(\cdot || \cdot) \geq 0$ (Gibbs' inequality), the global minimum of the objective is achieved exactly when the argument in Eq.5 becomes 0. To relate this to the sampling process, we note the invariance of \mathbb{KL} under diffeomorphic transformations, from which it follows that

$$\mathbb{KL}(p_{\hat{\phi}}(\mathbf{z} | \mathbf{x}) || p(\mathbf{z})) = 0. \quad (16)$$

Considering $p(\mathbf{z}) \perp p(\mathbf{x})$ and $p_{\hat{\phi}}(\mathbf{z} | \mathbf{x}) = p(\mathbf{z})$ (from Eq.16), this also implies $p_{\hat{\phi}}(\mathbf{z} | \mathbf{x}) \perp p(\mathbf{x})$, which means the latent output distribution is the same for any fixed \mathbf{x} we choose. This motivates the validity of taking samples from $p(\mathbf{z})$ and transforming them back using the condition, to generate samples from the posterior. By definition of the model, the generated samples $f_{\hat{\phi}}^{-1}(\mathbf{z}, \mathbf{x})$ with $\mathbf{z} \sim p(\mathbf{z})$ follow $p_{\hat{\phi}}(\boldsymbol{\theta} | \mathbf{x})$. The proposition therefore holds when the argument in Eq.5 is zero. \square

We now generalize our formulation to datasets with arbitrary numbers of observations. If we let the number of observations N vary and train a summary network $\tilde{\mathbf{x}} = h_{\psi}(\mathbf{x}_{1:N})$ together with the cINN, our main objective changes to:

$$\hat{\phi}, \hat{\psi} = \underset{\phi, \psi}{\operatorname{argmax}} \mathbb{E}_{p(\mathbf{x}, \boldsymbol{\theta}, N)} [\log p_{\phi}(\boldsymbol{\theta} | h_{\psi}(\mathbf{x}_{1:N}))] \quad (17)$$

and its Monte Carlo estimate to:

$$\hat{\phi}, \hat{\psi} = \underset{\phi, \psi}{\operatorname{argmin}} \frac{1}{M} \sum_{m=1}^M \left(\frac{\|f_{\phi}(\boldsymbol{\theta}^{(m)}; h_{\psi}(\mathbf{x}_{1:N}^{(m)}))\|_2^2}{2} - \log |\det (\mathbf{J}_{f_{\phi}}^{(m)})| \right) \quad (18)$$

In order to make the estimation of $p(\boldsymbol{\theta} | \mathbf{x}_{1:N})$ tractable, we assume that there exists a vector $\boldsymbol{\eta}$ of sufficient statistics that captures all information about $\boldsymbol{\theta}$ contained in $\mathbf{x}_{1:N}$ in a fixed-size (vector) representation. For $h_{\psi}(\mathbf{x}_{1:N})$ to be a useful estimator for $\boldsymbol{\eta}$, both should convey the same information about $\boldsymbol{\theta}$, as measured by the mutual information:

$$MI(\boldsymbol{\theta}, h_{\psi}(\mathbf{x}_{1:N})) \approx MI(\boldsymbol{\theta}, \boldsymbol{\eta}) \quad (19)$$

Since we do not know $\boldsymbol{\eta}$, we can enforce this requirement only indirectly by minimizing the Monte Carlo estimate of Eq.17. The following proposition states that, under perfect convergence, samples from a cINN still follow the true posterior given the outputs of a summary networks.

Proposition 2. Assume that we have a perfectly converged cINN f_ϕ and a perfectly converged summary network h_ψ . Assume also, that there exists a vector $\boldsymbol{\eta}$ of sufficient summary statistics for $\mathbf{x}_{1:N}$. Then, independently sampling $\mathbf{z} \sim p(\mathbf{z})$ and applying $f_\phi^{-1}(\mathbf{z}; h_\psi(\mathbf{x}_{1:N}))$ to each \mathbf{z} yields independent samples from $p(\boldsymbol{\theta} | \mathbf{x}_{1:N})$.

Proof. Perfect convergence of the networks under Eq.17 implies $\mathbb{KL}(p(\boldsymbol{\theta} | \mathbf{x}_{1:N}) || p_\phi(\boldsymbol{\theta} | h_\psi(\mathbf{x}_{1:N}))) = 0$. This, in turn, implies that $MI(\boldsymbol{\theta}, h_\psi(\mathbf{x}_{1:N})) = MI(\boldsymbol{\theta}, \boldsymbol{\eta})$, because a perfect match of the densities would be impossible if $h_\psi(\mathbf{x}_{1:N})$ contained less information about $\boldsymbol{\theta}$ than $\boldsymbol{\eta}$. Therefore, the proof reduces to that of **Proposition 1**. Note, that whenever the KL divergence is driven to a minimum, $h_\psi(\mathbf{x}_{1:N})$ is a *maximally informative statistic* [11]. \square

In summary, the approximate posteriors obtained by the BayesFlow method are correct if the summary and invertible networks are perfectly converged. In practice, however, perfect convergence is unrealistic and there are three sources of error which can lead to incorrect posteriors. The first is the Monte Carlo error introduced by using simulations from $g(\boldsymbol{\theta}, \boldsymbol{\xi})$ to approximate the expectation in Eq.17. The second is due to a summary network which may not fully capture the relevant information in the data or when sufficient summary statistics do not exist. The third is due to an invertible network which does not accurately transform the true posterior into the prescribed Gaussian latent space. Even though we can mitigate the Monte Carlo error by running the simulator $g(\boldsymbol{\theta}, \boldsymbol{\xi})$ more often, the latter two can be harder to detect and alleviate in a principled way. Nevertheless, recent work on *probabilistic symmetry* [6] and *algorithmic alignment* [52] can provide some guidelines on how to choose the right summary network for a particular problem. Additionally, the depth as well as the building blocks (to be explained shortly) of the invertible chain can be tuned to increase the expressiveness of the learned transformation from $\boldsymbol{\theta}$ -space to \mathbf{z} -space. The benefits of neural network depth have been confirmed both in theory and practice [28, 4], so we expect better performance in complex settings with increasing network depth.

2.3 Composing Invertible Networks

The basic building block of our cINN is the affine coupling block (ACB) [13]. Each ACB consists of four separate fully connected neural networks denoted as $s_1(\cdot), s_2(\cdot), t_1(\cdot), t_2(\cdot)$. An ACB performs an invertible non-linear transformation, which means that in addition to a parametric mapping $f_\phi : \mathbb{R}^D \rightarrow \mathbb{R}^D$ it also learns the inverse mapping $f_\phi^{-1} : \mathbb{R}^D \rightarrow \mathbb{R}^D$ for free. Denoting the input vector of f_ϕ as \mathbf{u} and the output vector as \mathbf{v} , it follows that $f_\phi(\mathbf{u}) = \mathbf{v}$ and $f_\phi^{-1}(\mathbf{v}) = \mathbf{u}$. Invertibility is achieved by splitting the input vector into two parts $\mathbf{u} = (\mathbf{u}_1, \mathbf{u}_2)$ with $\mathbf{u}_1 = \mathbf{u}_{1:D/2}$ and $\mathbf{u}_2 = \mathbf{u}_{D/2+1:D}$ (where $D/2$ is understood as a floor division) and performing the following operations on the split input:

$$\mathbf{v}_1 = \mathbf{u}_1 \odot \exp(s_1(\mathbf{u}_2)) + t_1(\mathbf{u}_2) \quad (20)$$

$$\mathbf{v}_2 = \mathbf{u}_2 \odot \exp(s_2(\mathbf{v}_1)) + t_2(\mathbf{v}_1) \quad (21)$$

where \odot denotes element-wise multiplication. The outputs $\mathbf{v} = (\mathbf{v}_1, \mathbf{v}_2)$ are then concatenated again and passed to the next ACB. The inverse operation is given by:

$$\mathbf{u}_2 = (\mathbf{v}_2 - t_2(\mathbf{v}_1)) \odot \exp(-s_2(\mathbf{v}_1)) \quad (22)$$

$$\mathbf{u}_1 = (\mathbf{v}_1 - t_1(\mathbf{u}_2)) \odot \exp(-s_1(\mathbf{u}_2)) \quad (23)$$

This formulation ensures that the Jacobian of the affine transformation is a strictly upper or a lower triangular matrix and therefore its determinant is very cheap to compute. Furthermore, the internal functions $s_1(\cdot), s_2(\cdot), t_1(\cdot), t_2(\cdot)$ can be represented by arbitrarily complex neural networks, which themselves need not be invertible, since they are only ever evaluated in the forward direction during both the forward and the inverse pass through the ACBs. In our applications, we parameterize the internal functions as fully connected neural networks with exponential linear units (ELU).

In order to ensure that the neural network architecture is expressive enough to represent complex distributions, we chain multiple ACBs, so that the output of each ACB becomes the input to the next one. In this way, the whole chain remains invertible from the first input to the last output and can be viewed as a single function parameterized by trainable parameters ϕ .

In our applications, the input to the first ACB is the parameter vector $\boldsymbol{\theta}$, and the output of the final ACB is a d -dimensional vector \mathbf{z} representing the non-linear transformation of the parameters. As described in the previous section, we ensure that \mathbf{z} follows a unit Gaussian distribution via optimization, that is, $p(\mathbf{z}) = \mathcal{N}_D(\mathbf{z} | \mathbf{0}, \mathbb{I})$. Fixed permutation matrices are used before each ACB to ensure that each axis of the transformed parameter space \mathbf{z} encodes information from all components of $\boldsymbol{\theta}$.

In order to account for the observed data, we feed the learned summary vectors into all internal networks of each ACB (explained shortly). Intuitively, in this way we realize the following process: the forward pass maps data-generating

parameters θ to z -space using conditional information from the data $\mathbf{x}_{1:N}$, while the inverse pass maps data points from z -space to the data-generating parameters of interest using the same conditional information.

2.4 Summary Network

Since the number of observations usually varies in practical scenarios (e.g., different number of measurements or time points) and since datasets might exhibit various redundancies, the cINN can profit from some form of dimensionality reduction. As previously mentioned, we want to avoid information loss through restrictive hand-crafted summary statistics and, instead, learn the most informative summary statistics directly from data. Therefore, instead of feeding the raw simulated or observed data to each ACB, we pass the data through an additional summary network to obtain a fixed-sized vector of learned summary statistics $\tilde{\mathbf{x}} = h_{\psi}(\mathbf{x}_{1:N})$.

The architecture of the summary network should be aligned with the probabilistic symmetry of the observed data. An obvious choice for time series-data is an LSTM-network [15], since recurrent networks can naturally deal with long sequences of variable size. Another choice might be a 1D fully convolutional network [29], which has already been applied in the context of likelihood-free inference [41]. A different architecture is needed when dealing with *i.i.d.* samples of variable size. Such data are often referred to as *exchangeable*, or *permutation invariant*, since changing the order of individual elements does not change the associated likelihood or posterior. In other words, if $\mathbb{S}_N(\cdot)$ is an arbitrary permutation of N elements, the following should hold for the posterior:

$$p(\theta | \mathbf{x}_{1:N}) = p(\theta | \mathbb{S}_N(\mathbf{x}_{1:N})) \quad (24)$$

Following [6], we encode probabilistic permutation invariance by implementing a permutation invariant function through an equivariant non-linear transformation followed by a pooling operator (e.g., sum or mean) and another non-linear transformation:

$$\tilde{\mathbf{x}} = h_{\psi_1} \left(\sum_{i=1}^N h_{\psi_2}(\mathbf{x}_i) \right) \quad (25)$$

where h_{ψ_1} and h_{ψ_2} are two different fully connected neural networks. In practice, we stack multiple equivariant and invariant functions into an invariant network in order to achieve higher expressiveness [6].

We optimize the parameters ψ of the summary network jointly with those of the cINN chain via backpropagation. Thus, training remains completely end-to-end, and BayesFlow learns to generalize to datasets of different sizes by suitably varying N during training of a permutation invariant summary network or varying sequence length during training of a recurrent/convolutional network.

To incorporate the observed or simulated data $\mathbf{x}_{1:N}$, each of the internal networks of each ACB is augmented to take the learned summary vector $\tilde{\mathbf{x}}$ of the data as an additional input. The output of each ACB then becomes:

$$\mathbf{v}_1 = \mathbf{u}_1 \odot \exp(s_1(\mathbf{u}_2, \tilde{\mathbf{x}})) + t_1(\mathbf{u}_2, \tilde{\mathbf{x}}) \quad (26)$$

$$\mathbf{v}_2 = \mathbf{u}_2 \odot \exp(s_2(\mathbf{v}_1, \tilde{\mathbf{x}})) + t_2(\mathbf{v}_1, \tilde{\mathbf{x}}) \quad (27)$$

Thus, a complete pass through the entire conditional invertible chain can be expressed as $f_{\phi}(\theta; \tilde{\mathbf{x}}) = \mathbf{z}$ together with the inverse operation $f_{\phi}^{-1}(\mathbf{z}; \tilde{\mathbf{x}}) = \theta$. The inverse transformation during inference is depicted in Figure 2.

2.5 Putting It All Together

Algorithm 1 describes the essential steps of the BayesFlow method using an arbitrary summary network and employing an online learning approach.

The backpropagation algorithm works by computing the gradients of the loss function with respect to the parameters of the neural networks and then adjusting the parameters, so as to drive the loss function to a minimum. We experienced no instability or convergence issues during training with the loss function given by Eq.18. Note, that steps 3-14 and 18-22 of **Algorithm 1** can be executed in parallel with GPU support in order to dramatically accelerate convergence and inference. Moreover, steps 18-22 can be applied in parallel to an arbitrary number of observed datasets after convergence of the networks (see Figure 2 for a full graphical illustration).

In what follows, we apply BayesFlow to two toy models with a unimodal and multimodal posteriors, respectively, and then use it to perform Bayesian inference on challenging models from population dynamics, cognitive science, epidemiology, and ecology.¹ We deem these models suitable for an initial validation, since they differ widely in the

¹Code and simulation scripts for all current applications are available at <https://github.com/stefanradev93/cINN>.

Algorithm 1 Amortized Bayesian inference with the BayesFlow method

```

1: Training phase (online learning with batch size  $M$ ):
2: repeat
3:   Sample number of observations  $N \sim \mathcal{U}(N_{min}, N_{max})$ .
4:   for  $m = 1, \dots, M$  do
5:     Sample model parameters from prior:  $\theta^{(m)} \sim p(\theta)$ .
6:     for  $i = 1, \dots, N$  do
7:       Sample a noise instance:  $\xi_i \sim p(\xi)$ .
8:       Run the simulation (cf. Eq.1) to create a synthetic observation:  $x_i^{(m)} = g(\theta^{(m)}, \xi_i)$ .
9:     end for
10:    Pass the dataset  $x_{1:N}^{(m)}$  through the summary network:  $\tilde{x}^{(m)} = h_\psi(x_{1:N}^{(m)})$ .
11:    Pass  $(\theta^{(m)}, \tilde{x}^{(m)})$  through the inference network in forward direction:  $z^{(m)} = f_\phi(\theta^{(m)}; \tilde{x}^{(m)})$ .
12:  end for
13:  Compute loss according to Eq.18 from the training batch  $\{(\theta^{(m)}, \tilde{x}^{(m)}, z^{(m)})\}_{m=1}^M$ .
14:  Update neural network parameters  $\phi, \psi$  via backpropagation.
15: until convergence to  $\hat{\phi}, \hat{\psi}$ 
16:
17: Inference phase (given observed or test data  $x_{1:N}^o$ ):
18: Pass the observed dataset through the summary network:  $\tilde{x}^o = h_{\hat{\psi}}(x_{1:N}^o)$ .
19: for  $l = 1, \dots, L$  do
20:   Sample a latent variable instance:  $z^{(l)} \sim \mathcal{N}_D(\mathbf{0}, \mathbb{I})$ .
21:   Pass  $(\tilde{x}^o, z^{(l)})$  through the inference network in inverse direction:  $\theta^{(l)} = f_{\hat{\phi}}^{-1}(z^{(l)}; \tilde{x}^o)$ .
22: end for
23: Return  $\{\theta^{(l)}\}_{l=1}^L$  as a sample from  $p(\theta | x_{1:N}^o)$ 

```

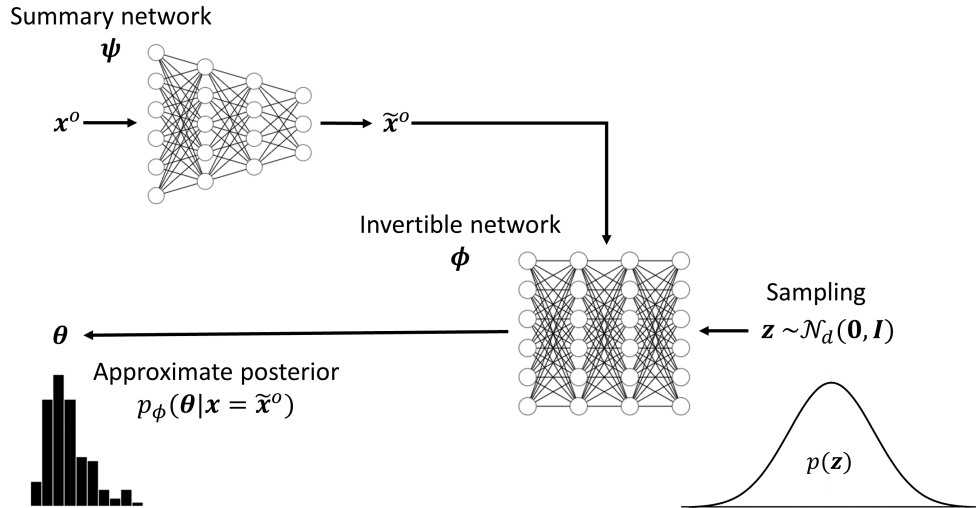


Figure 2: Inference with pre-trained summary and inference networks. The posterior is approximated given real observed data via independent samples from a learned pushforward distribution. Thus, knowledge about the mapping between data and parameters (the inverse model) is compactly encoded within the weights of the two networks.

generative mechanisms they implement and the observed data they aim to explain. Therefore, good performance on these disparate examples underlines the broad empirical utility of the BayesFlow method. Details for models' setup can be found in **Appendix C**.

3 Experiments

3.1 Training the Networks

We train all invertible and summary networks described in this paper jointly via backpropagation. For all following experiments, we use the Adam optimizer with a starter learning rate of 10^{-3} and an exponential decay rate of .95. We perform 50 000 to 100 000 iterations (i.e., mini-batch update steps) for each experiment, and report the results obtained by the converged networks. Note, that we did not perform an extensive search for optimal values of network hyperparameters, but use a default BayesFlow with 5 to 10 ACBs and a summary vector of size 128 for all examples in this paper (see **Appendix C** for more details on summary network architectures). All networks were implemented in Python using the *TensorFlow* library [1] and trained on a single-GPU machine equipped with NVIDIA® GTX1060 graphics card.

Regarding the data generation step, we take an approach which incorporates ideas from *online learning* [36] where data are simulated by Eq.1 on demand. Correspondingly, a dataset $\mathbf{x}_{1:N}$, or a batch of M datasets $\{\mathbf{x}_{1:N}^{(m)}\}_{m=1}^M$, is generated on the fly and then passed through the neural network. This training approach has the advantage that the network never *experiences* the same input data twice. Moreover, training can continue as long as the network keeps improving (i.e., the loss keeps decreasing), since overfitting in the classical sense is nearly impossible. However, if simulations are computationally expensive and researchers need to experiment with different networks or training hyperparameters, it might be beneficial to store and re-use simulations, since simulation and training in online learning are tightly intertwined.

Once the networks have converged, we store the trained networks and use them to perform amortized inference on a separate validation set of datasets. The pre-trained networks can also be shared among a research community so that multiple researchers/labs can benefit from the amortization of inference.

3.2 Performance Validation

To evaluate the performance of BayesFlow in the following application examples, we consider a number of different metrics:

- Normalized root mean squared error (NRMSE) - to assess accuracy of point-estimates in recovering ground-truth parameter values;
- Coefficient of determination (R^2) - to assess the proportion of variance in ground-truth parameters that is captured by the point estimates;
- Re-simulation error (Err_{sim}) - to assess the predictive mismatch between the true data distribution and the data distribution generated with the estimated parameters (i.e., posterior predictive check);
- Calibration error (Err_{cal} , [2]) - to assess the coverage of the approximate posteriors (i.e., whether credibility intervals are indeed credible);
- Simulation-based calibration (SBC, [46]) - to visually detect systematic biases in the approximate posteriors;

Details for computing all metrics are given in **Appendix B**.

3.3 Proof of Concept: Multivariate Normal Distribution

As a proof-of-concept, we apply the BayesFlow method to recover the posterior mean vector of a toy multivariate normal (MVN) example. For a single D -dimensional MVN vector, the forward model is given by:

$$\boldsymbol{\mu}^{(m)} \sim \mathcal{N}_D(\boldsymbol{\mu} | \mathbf{0}, \mathbb{I}) \quad (28)$$

$$\mathbf{x}^{(m)} \sim \mathcal{N}_D(\mathbf{x} | \boldsymbol{\mu}^{(m)}, \boldsymbol{\Sigma}) \quad (29)$$

where in this illustrative case we assume a single D -dimensional sample per observation ($N = 1$). If the covariance matrix $\boldsymbol{\Sigma}$ is known, the posterior of the mean vector $\boldsymbol{\mu}$ has a closed-form which is also a MVN $p(\boldsymbol{\mu} | \mathbf{x}, \boldsymbol{\Sigma}) = \mathcal{N}_d(\boldsymbol{\mu} | \mathbf{m}, \boldsymbol{\Lambda})$ with posterior precision matrix given by $\boldsymbol{\Lambda}^{-1} = \mathbb{I} + \boldsymbol{\Sigma}^{-1}$ and posterior mean given by $\mathbf{m} = \boldsymbol{\Lambda} \boldsymbol{\Sigma}^{-1} \mathbf{x}$

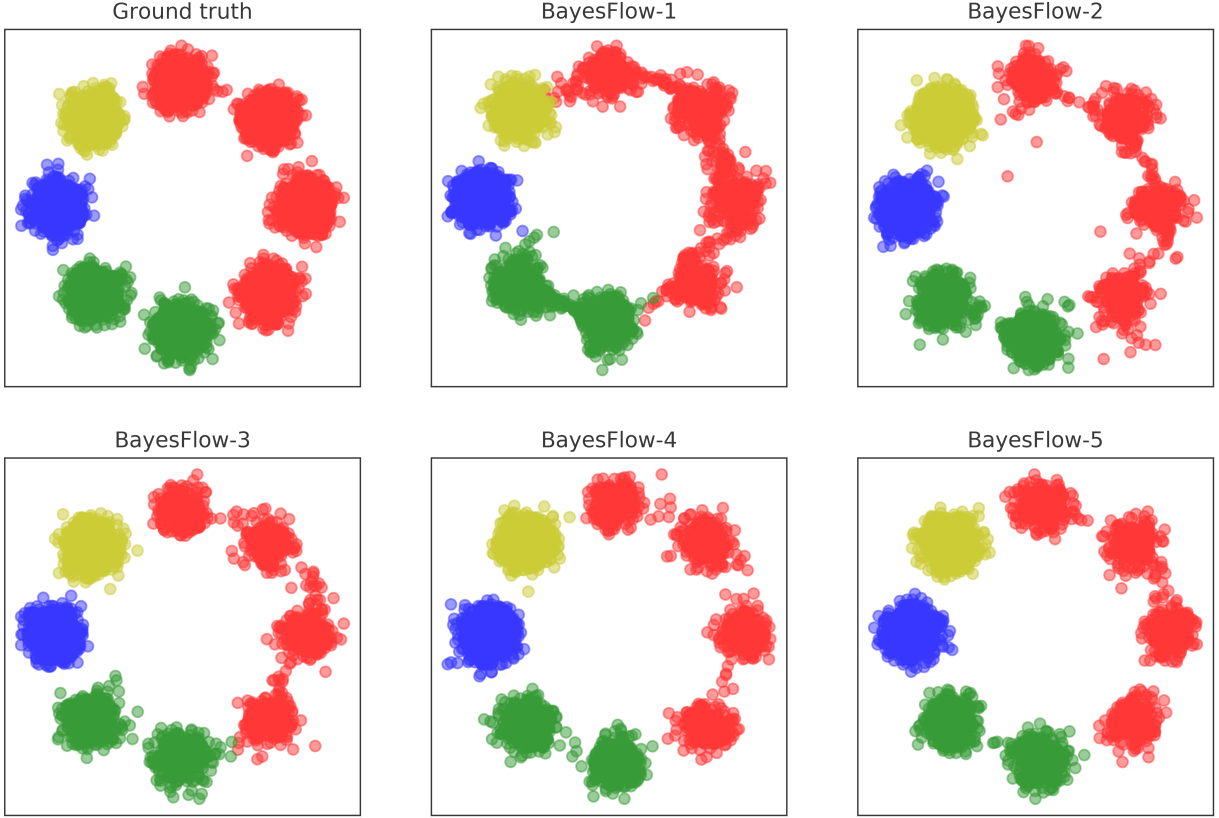


Figure 3: Results on the GMM toy example with colors indicating cluster assignments. Approximation of the multimodal posterior become closer to the ground truth distribution with increasing depth (number of ACBs) of the conditional invertible network.

[7]. We can thus generate multiple batches of the form $\{(\mathbf{x}^{(m)}, \boldsymbol{\mu}^{(m)})\}_{m=1}^M$ and pass them directly through an invertible network. Since the ground-truth posterior is Gaussian, we can compute the KL divergence as a measure of mismatch between the true and approximate posteriors in closed form.

We run three experiments with $D \in \{5, 50, 500\}$ where the size of the ACB blocks was doubled for each successive D . To assess results, we compute the R^2 and NRMSE between approximate and true means as well as the KL divergence between approximate and true distributions on 100 test datasets. To compute the approximate covariance matrix, we draw 5000 samples from the approximate posteriors for $D = 5$ and $D = 50$ and 50000 samples for $D = 500$.

The KL divergence for the 5-D and 50-D MVNs reached essentially 0 after 2-3 epochs of 1000 iterations indicating that this is an easy problem for BayesFlow, and almost perfect recovery of the true posteriors is possible. The KL divergence for the 500-D MVN model reached 0.37 after 50 epochs, which represents a negligible increase in entropy relative to the true posterior (0.05% nats) and indicates decent approximation in light of the high dimensionality of the problem.

3.4 Multimodal Posterior - Gaussian Mixture Model

In order to test whether the BayesFlow method can recover multimodal posteriors, we apply it to a generative Gaussian mixture model (GMM). Multimodal posteriors arise in practice, for example, when forward models are defined as mixtures between different processes, or when models exhibit large multivariate trade-offs in their parameter space (e.g., there are multiple separate regions of posterior density with plausible parameter values). Therefore, it is important to show that our method is able to capture such behavior and does not suffer from mode collapse.

Following [2], we construct a scenario in which the observed data \mathbf{x} is a one-hot encoded vector representing one of the *hard* labels *red*, *green*, *blue*, or *yellow* (i.e., a single observation, thus $N = 1$). The parameters $\boldsymbol{\theta} = (\theta_1, \theta_2)$ are the 2D

coordinates² of points drawn from a mixture of eight Gaussian clusters with centers distributed around the origin in a clockwise manner and unit variance (see Figure 3, upper left). The first four clusters are assigned the label *red*, the next two the label *green*, and the remaining two the labels *blue* and *yellow*. The posterior $p(\theta | x)$ is composed of the clusters indexed by the corresponding label. We perform the experiment multiple times by increasing the depth of the BayesFlow starting from 1 ACB block up to 5 ACB blocks. In this way, we can investigate the effects of cINN depth on the quality of the approximate multimodal posteriors. We train each BayesFlow for 50 epochs and draw 8000 samples from the approximate posteriors obtained by the trained models.

Results for all BayesFlows are depicted in Figure 3. We observe that approximations profit from having a deeper cINN chain, with cluster separation becoming clearer when using more ACBs. This confirms that our method is capable of recovering multimodal posteriors.

3.5 Stochastic Time-Series Model - The Ricker Model

In the following, we estimate the parameters of a well-known discrete stochastic population dynamics model [51]. With this example, we are pursuing several goals: First, we want to demonstrate that the BayesFlow method is able to accurately recover the parameters of an actual model with intractable likelihood by learning summary statistics from raw data. Second, we show that BayesFlow can deal adequately with parameters that are completely unrelated to the data by reducing estimates to the corresponding parameters’ prior. Third, we compare the global performance of the BayesFlow method to that of related methods capable of amortized likelihood-free inference. Finally, we demonstrate the desired posterior contraction and improvement in estimation with increasing number of observations.

Discrete population dynamics models describe how the number of individuals in a population changes over discrete units of time [51]. In particular, the Ricker model describes the number of individuals x_t in generation t as a function of the expected number of individuals in the previous generation by the following non-linear equations:

$$x_t \sim \text{Pois}(\rho N_t) \quad (30)$$

$$\xi_t \sim \mathcal{N}(0, \sigma^2) \quad (31)$$

$$N_{t+1} = r N_t e^{-N_t + \xi_t} \quad (32)$$

for $t = 1, \dots, T$ where N_t is the expected number of individuals at time t , r is the growth rate, ρ is a scaling parameter and ξ_t is random Gaussian noise. The likelihood function for the Ricker model is not available in closed form, and the model is known to exhibit chaotic behavior [33]. Thus, it is a suitable candidate for likelihood-free inference. The parameter estimation task consists of recovering $\theta = (\rho, r, \sigma)$ from the observed one-dimensional time-series data $x_{1:T}$ where each $x_t \in \mathbb{N}$.

What if the data does not contain any information about a particular parameter? In this case, any good estimation method should detect this, and return the prior of the particular parameter. To test this, we append a random uniform variable $u \sim \mathcal{U}(0, 1)$ to the parameter vector θ and train BayesFlow with this additional dummy parameter. We expect that the networks ignore this dummy parameter, that is, we assume that the estimated posterior of u resembles the uniform prior.

We compare the performance of BayesFlow to the following recent methods capable of amortized likelihood-free inference: conditional variational autoencoder (cVAE) [35], cVAE with autoregressive flow (cVAE-IAF) [26], *DeepInference* with heteroscedastic loss [41], approximate Bayesian computation with an LSTM neural network for learning informative summary statistics (ABC-NN) [22] and quantile random forest (ABC-RF) [43]. For training the models, we simulate time-series from the Ricker model with varying lengths. The number of time points T is drawn from a uniform distribution $T \sim \mathcal{U}(100, 500)$ at each training iteration.

All neural network methods were trained for 100 epochs with 1000 iterations each on simulated data from the Ricker model. The ABC-RF method was fitted on a reference table with 200 000 datasets, since the method does not allow for online learning and increasing the reference table did not seem to improve performance. In order to avoid using hand-crafted summary statistics for the ABC-RF method, we input summary vectors obtained by applying the summary network trained jointly with the cINN. Thus, the ABC-RF method has the advantage of using maximally informative statistics as input. We validate the performance of all methods on an independent test set of 500 datasets generated with $T = 500$. We report performance metrics for each method and each parameter in Table 1.

Parameters r and ρ seem to be well recoverable by all methods considered here. The σ parameter turns out to be harder to estimate, with BayesFlow and the ABC-NN method performing best. Further, BayesFlow performs very

²Note that this is not the typical GMM setup, as we construct the example such that the mixture assignments (labels) are observed and the data coordinates are the latent parameters.

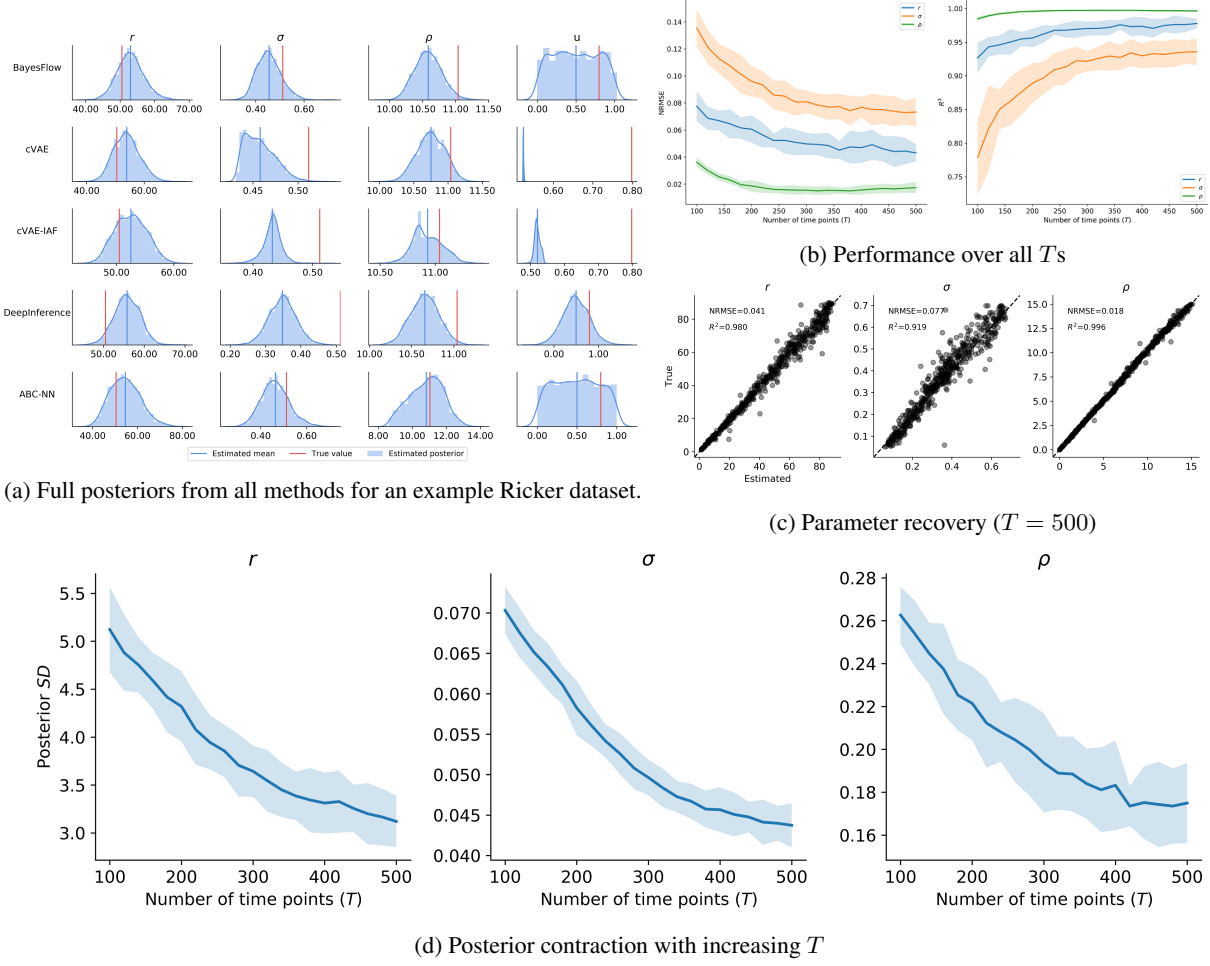


Figure 4: Results on the Ricker model. **(a)** Approximate posteriors obtained by all implemented methods on a single Ricker dataset. Note that only BayesFlow and ABC-NN are able to approximate the uniform posterior of u ; **(b)** NRMSE and R^2 performance metrics over all T s obtained by the BayesFlow method. We observe that parameter estimation remains good over all T s, and becomes progressively better as more data is available (shaded regions indicate bootstrap 95% CIs); **(c)** Parameter recovery with BayesFlow for the maximum number of generations used during training ($T = 500$); **(d)** Posterior contraction in terms of posterior standard deviation for each parameter across increasing number of available generations (shaded regions indicate bootstrap 95% CIs).

well across all parameters and metrics. Importantly, the calibration error Err_{cal} obtained by BayesFlow is always low, indicating that the shape of the approximate posterior closely matches that of the true posteriors. Variational methods (cVAE, cVAE-IAF) experience some problems recovering the posterior of σ . The ABC-NN and ABC-RF methods seem to recover point estimates with high accuracy but the approximate posteriors of the former exhibit relatively high calibration error. The ABC-RF method can only estimate posterior quantiles, so no comparable calibration metric could be computed.

Further results are depicted in Figure 4. Inspecting the full posteriors obtained by all methods on an example test dataset, we note that only BayesFlow and the ABC-NN methods are able to recover the uninformative posterior distribution of the dummy noise variable u (Figure 4a). Moreover, the importance of a Bayesian treatment of the Ricker model becomes clear when looking at the posteriors of σ . On most test datasets, the posterior density spreads over the entire prior range (high posterior variance) indicating large uncertainty in the obtained estimates. Moreover, the shapes of the marginal parameter posteriors vary widely across validation datasets, which highlights the importance of avoiding *ad hoc* restrictions on allowed posterior shapes (see Figure S5 for examples). We also observe that parameter estimation with BayesFlow becomes increasingly accurate when more time points are available (Figure 4b). Parameter recovery is especially good with the maximum number of time points (see Figure 4c). Finally, (Figure 4d) reveals a notable posterior contraction across increasing number of time points available to the summary network.

Table 1: Performance results on the Ricker model across all estimation methods

		BayesFlow	cVAE	cVAE-IAF	DeepInference	ABC-NN	ABC-RF
Err_{cal}	r	0.017 ± 0.007	0.014 ± 0.007	0.058 ± 0.017	0.122 ± 0.016	0.164 ± 0.015	-
	σ	0.013 ± 0.007	0.419 ± 0.011	0.382 ± 0.013	0.184 ± 0.021	0.119 ± 0.014	-
	ρ	0.084 ± 0.018	0.121 ± 0.017	0.188 ± 0.018	0.111 ± 0.019	0.283 ± 0.012	-
$NRMSE$	r	0.041 ± 0.002	0.047 ± 0.004	0.047 ± 0.006	0.052 ± 0.003	0.053 ± 0.003	0.044 ± 0.004
	σ	0.077 ± 0.005	0.137 ± 0.004	0.124 ± 0.006	0.108 ± 0.004	0.077 ± 0.004	0.081 ± 0.005
	ρ	0.018 ± 0.001	0.016 ± 0.002	0.019 ± 0.002	0.019 ± 0.002	0.033 ± 0.002	0.021 ± 0.001
R^2	r	0.980 ± 0.003	0.973 ± 0.005	0.973 ± 0.007	0.968 ± 0.005	0.966 ± 0.004	0.977 ± 0.004
	σ	0.919 ± 0.011	0.745 ± 0.020	0.792 ± 0.020	0.841 ± 0.014	0.919 ± 0.010	0.912 ± 0.011
	ρ	0.996 ± 0.001	0.997 ± 0.001	0.996 ± 0.001	0.996 ± 0.001	0.986 ± 0.002	0.994 ± 0.001
Err_{sim}	-	0.038 ± 0.001	0.041 ± 0.001	0.042 ± 0.001	0.041 ± 0.001	0.048 ± 0.002	0.041 ± 0.002

Note: For each parameter, bootstrapped means (± 1 standard error) of different performance metrics are displayed for all tested methods. For each metric and each parameter, the best performance across methods is printed in bold font.

3.6 A Model of Perceptual Decision Making - The Lévy-Flight Model

In the following, we estimate the parameters of a stochastic differential equation model of human decision making. We perform the first Bayesian treatment of the recently proposed Lévy-Flight Model (LFM), as its intractability has so far rendered traditional non-amortized Bayesian inference methods prohibitively slow [49].

With this example, we first want to show empirically that BayesFlow is able to deal with *i.i.d.* datasets of variable size arising from N independent runs of a complex stochastic simulator. For this, we inspect global performance of BayesFlow over a wide range of dataset sizes. Additionally, we want to show the advantage of amortized inference compared to case-based inference in terms of efficiency and recovery. For this, we apply BayesFlow along with four other recent methods for likelihood-free inference to a single dataset and show that in some cases the speed advantage of amortized inference becomes noticeable even after as few as 5 datasets. Crucially, researchers often fit the same models to different datasets, so if a pre-trained model exists, it would present a huge advantage in terms of efficiency and productivity.

We focus on the family of evidence accumulator models (EAMs) which describe human decision making by a set of neurocognitively motivated parameters [42]. EAMs are most often applied to choice reaction times (RT) data to obtain an estimate of the underlying processes governing categorization and (perceptual) decision making. In its most general formulation, the forward model of EAMs takes the form of a stochastic ordinary differential equation (ODE):

$$dx = vdt + \xi\sqrt{dt} \quad (33)$$

where dx denotes a change in activation of an accumulator, v denotes the average speed of information accumulation (often termed the drift rate), and ξ represents a stochastic additive component, usually modeled as coming from a Gaussian distribution centered around 0: $\xi \sim \mathcal{N}(0, c^2)$.

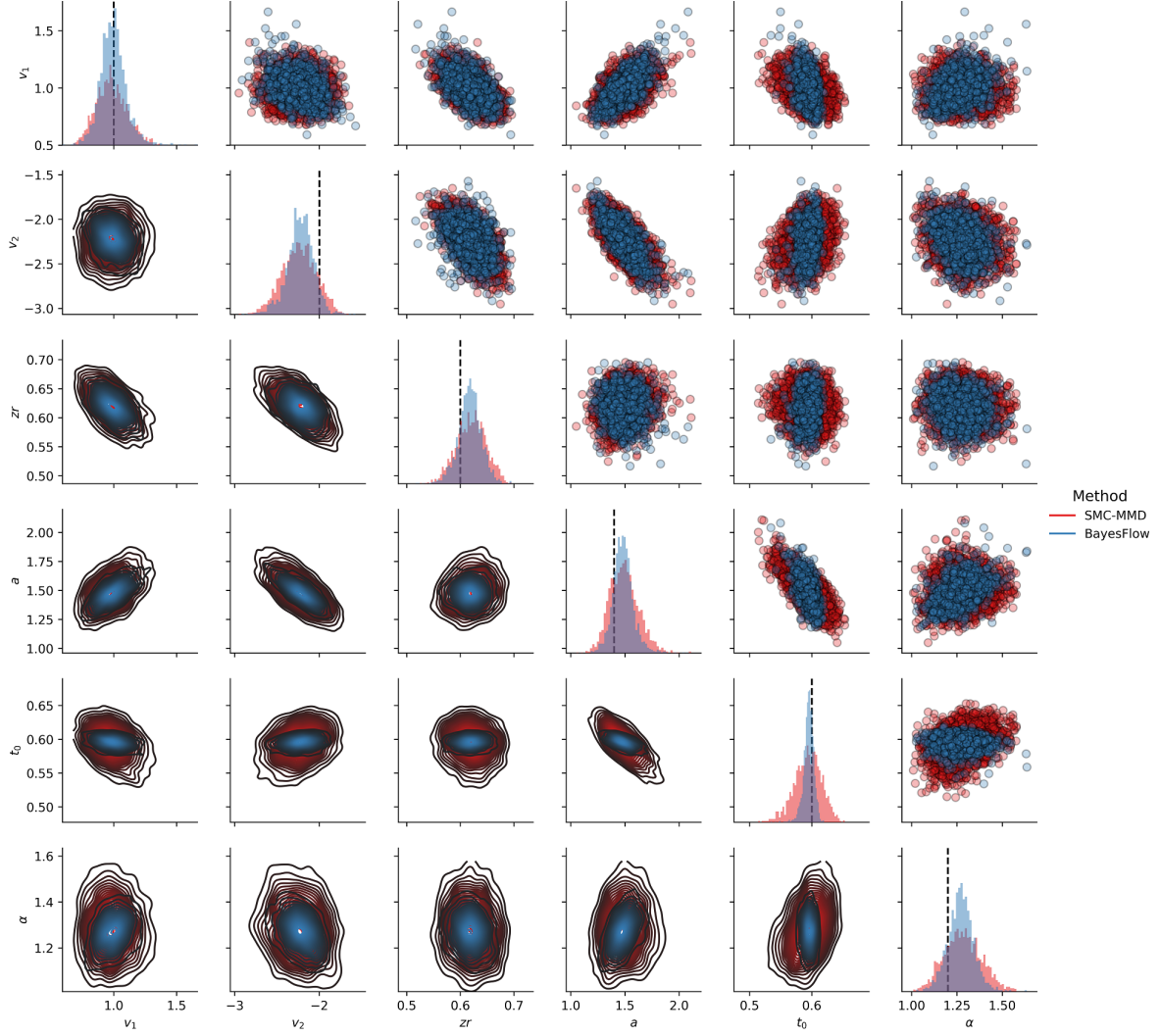
EAMs are particularly amenable for likelihood-free inference, since the likelihood of most interesting members of this model family turn out to be intractable [34]. This intractability has precluded many interesting applications and empirically driven model refinements. Here, we apply BayesFlow to estimate the parameters of the recently proposed Lévy-Flight Model (LFM) [49]. The LFM assumes an α -stable noise distribution of the evidence accumulation process which allows to model discontinuities in the decision process. However, the inclusion of α -stable noise (instead of the typically assumed Gaussian noise) leads to a model with intractable likelihood:

$$dx = vdt + \xi dt^{1/\alpha} \quad (34)$$

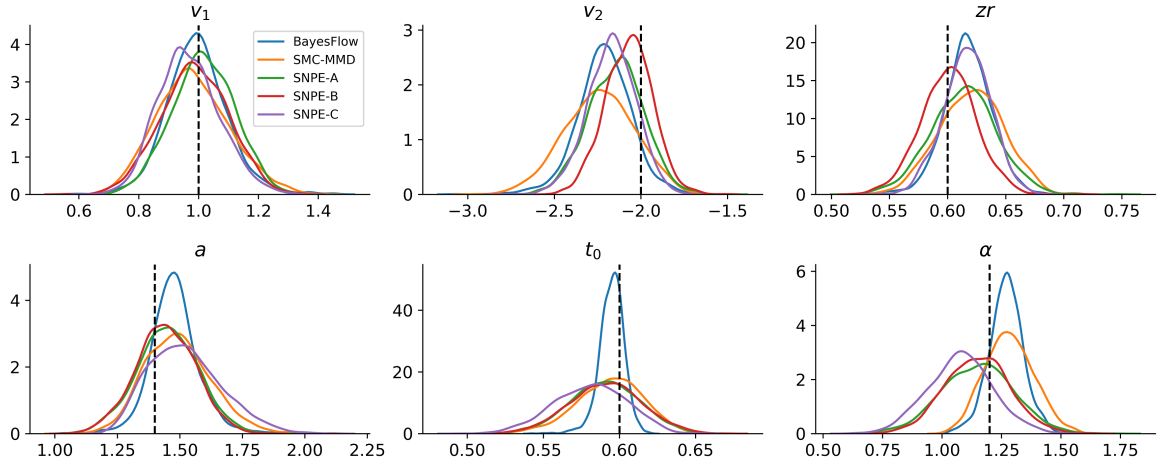
$$\xi \sim \text{AlphaStable}(\alpha, 0, 1, 0) \quad (35)$$

where α controls the probability of outliers in the noise distribution. The LFM has three additional parameters: the threshold a determining the amount of evidence needed for the termination of a decision process; a relative starting point, zr , determining the amount of starting evidence available to the accumulator before the actual decision alternatives are presented; and an additive non-decision time t_0 .

During training of the networks, we simulate response times data from two experimental conditions with two different drift rates, since such a design is often encountered in psychological research. The parameter estimation task is thus to recover the parameters $\theta = (v_0, v_1, a, t_0, zr, \alpha)$ from two-dimensional *i.i.d.* RT data $\mathbf{x}_{1:N}$ where each $\mathbf{x}_i \in \mathbb{R}^2$ represents RTs obtained in the two conditions. The number of trials is drawn from a uniform distribution



(a) Joint posteriors from BayesFlow and SMC-MMD



(b) Marginal posteriors from all methods

Figure 5: Comparison results on the LFM model. **(a)** Marginal and bivariate posteriors obtained by BayesFlow and SMC-MMD on the single validation dataset. We observe markedly better sharpness in the BayesFlow posteriors; **(b)** Marginal posteriors obtained from all methods under comparison.

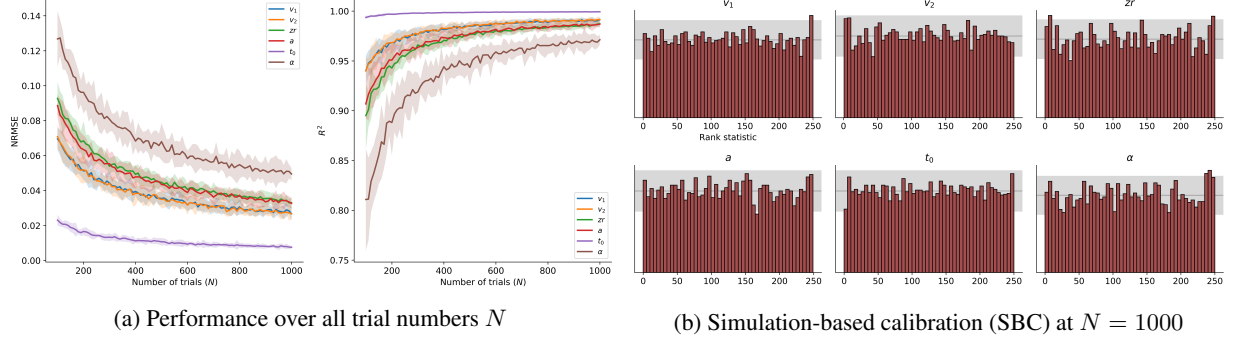


Figure 6: BayesFlow results obtained on the LFM model.

Table 2: Speed of inference and break-even for amortized inference for the LFM model

	Upfront Training	Inference (1 dataset)	Inference (500 datasets)	Break-even after
BayesFlow	23.2 h	60 ms	3.7 s	-
SMC-MMD	-	5.5 h	2700 h	5 datasets
SNPE-A	-	0.65 h	325 h	37 datasets
SNPE-B	-	0.65 h	325 h	37 datasets
SNPE-C	-	0.35 h	175 h	75 datasets

Note: Inference times for 500 datasets as well as the number of datasets for break-even with BayesFlow for the SMC-MMD, SNPE-A, SNPE-B, and SNPE-C methods are extrapolated from the wall-clock running time on a single dataset, so these are approximate quantities.

$N \sim \mathcal{U}(100, 1000)$ at each training iteration. Training the networks took a little less than a day with the online learning approach. Inference on 1000 datasets with 2000 posterior samples per parameter took approximately 7.39 seconds.

In order to investigate whether amortized inference is advantageous for this model, we additionally apply a version of the SMC-ABC algorithm available in the *pyABC* package [27] to a single dataset with $N = 500$. Since no sufficient summary statistics are available for EAM data, we apply the maximum mean discrepancy (MMD) metric as a distance between the full raw empirical RT distributions, in order to prevent information loss [39]. Since the MMD is expensive to compute, we use a GPU implementation to ensure that computation of MMD is not a bottleneck for the comparison. In order to achieve good approximation with 2000 samples from the SMC-MMD approximate posterior, we run the algorithm for 20 populations with a final rejection threshold $\epsilon = 0.04$. We also draw 2000 samples from the approximate posterior obtained by applying our pre-trained BayesFlow networks to the same dataset.

Along SMC-MMD, we apply three recent methods for neural density estimation, SNPE-A [37], SNPE-B [30], and SNPE-C ([16], also dubbed APT). Since these methods all depend on summary statistics of the data, we compute the first 6 moments of each empirical response time distribution as well as the fractions of correct/wrong responses. We train each method for a single round with 100 epochs and 5000 simulated datasets, in order to keep running time at a minimum. Also, we did not observe improvement in performance when training for more than one round. For each model, we sample 2000 samples from the approximate joint posterior to align the number of samples with those obtained via SMC-MMD.

The comparison results are depicted in Figure 5. We first focus on the comparison with SMC-MMD on the single dataset. Figure 5a depicts marginal and bivariate posteriors obtained by BayesFlow and SMC-MMD. The approximate posteriors of BayesFlow appear noticeably sharper. Observing the SCB plots (Figure 6b), we can conclude that the approximate posteriors of BayesFlow mirror the sharpness of the true posterior, since otherwise the SCB plots would show marked deviations from uniformity. Further, Figure 5b depicts the marginal posteriors obtained from the application of each method. Noticeably, performance and sharpness varies across the methods and parameters, with all methods yielding good point-estimate recovery via posterior means in terms of the NRMSE and R^2 metrics.

Importantly, Table 2 summarizes the advantage of amortized inference for the LFM model in terms of efficiency. For instance, compared to SMC-MMD, the extra effort of learning a global BayesFlow model upfront is worthwhile even after as few as 5 datasets, as inference with SMC-MMD would have taken more than a day to finish. On the other hand, the break-even for SNPE-C/APT occurs after 75 datasets, so in cases where only a few dozens of datasets are

considered, case-based inference might be preferable. However, the difficulties in manually finding meaningful and efficiently computable summary statistics may eat up possible savings even in this situation. We acknowledge that our choices in this respect might be sub-optimal, so performance comparisons should be treated with some caution.

We note, that after a day of training, the pre-trained networks of BayesFlow take less than 5 seconds to perform inference on 500 datasets even with maximum number of trials $N = 1000$. Using the case-based SMC-MMD algorithm, 500 inference runs would have taken more than half a year to complete. We also note, that parallelizing separate inference threads across multiple cores or across nodes of a (GPU) computing cluster can dramatically increase the wall-clock speed of the case-based methods considered here. However, the same applies to BayesFlow training, since its most expensive part, the simulation from the forward model, would profit the most from parallel computing.

The global performance of BayesFlow over all validation datasets and all trial sizes N is depicted in (Figure 6). First, we observe excellent recovery of all LFM parameters with NRMSEs ranging between 0.008 and 0.048 and R^2 between 0.972 and 0.99 for the maximum number of trials. Importantly, estimation remains very good across all trial numbers, and improves as more trials become available (Figure 6a). The parameter α appears to be most challenging to estimate, requiring more data for good estimation, whereas the non-decision time parameter t_0 can be recovered almost perfectly for all trial sizes. Last, the SCB histograms indicate no systematic deviations across the marginal posteriors (Figure 6b).

3.7 Stochastic Differential Equations - The SIR Epidemiology Model

With this example, we want to further corroborate the excellent global performance and probabilistic calibration observed for the LFM model on a non-*i.i.d.* stochastic ODE model. For this, we study a compartmental model from epidemiology, whose output comprises variable-sized multidimensional and inter-dependent time-series. It is therefore of interest to investigate how our method performs when applied to data which is the direct output of an ODE simulator.

Compartmental models in epidemiology describe the stochastic dynamics of infectious diseases as they spread over a population of individuals [23, 20]. The parameters of these models encode important characteristics of diseases, such as infection and recovery rates. The stochastic SIR model describes the transition of N individuals between three discrete states – susceptible (S), infected (I), and recovered (R) – whose dynamics follow the equations:

$$\Delta S = -\Delta N_{SI} \quad (36)$$

$$\Delta I = \Delta N_{SI} - \Delta N_{IR} \quad (37)$$

$$\Delta R = \Delta N_{IR} \quad (38)$$

$$\Delta N_{SI} \sim \text{Binomial}(S, 1 - \exp\left(-\beta \frac{I}{N} \Delta t\right)) \quad (39)$$

$$\Delta N_{IR} \sim \text{Binomial}(I, 1 - \exp(-\gamma \Delta t)) \quad (40)$$

where $S + I + R = N$ give the number of susceptible, infected, and recovered individuals, respectively. The parameter β controls the transition rate from being susceptible to infected, and γ controls the transition rate from being infected to recovered. The above listed stochastic system has no analytic solution and thus requires numerical simulation methods for recovering parameter values from data. Cast as a parameter estimation task, the challenge is to recover $\theta = \{\beta, \gamma\}$ from three dimensional time-series data $\mathbf{x}_{1:T}$ where each $\mathbf{x}_t \in \mathbb{N}^3$ is a triple containing the number of susceptible (S), number of infected (I), and recovered (R) individuals at time t .

During training of the networks, we simulate time-series from the stochastic SIR model with varying lengths. The number of time points T is drawn from a uniform distribution $T \sim \mathcal{U}(200, 500)$ at each training iteration. For small T , the system has not yet reached an equilibrium (i.e., not all individuals have transitioned from being I to R). It is especially interesting to see if BayesFlow can recover the rate parameters, while the process dynamics are still unfolding over time. Training the networks took approximately two hours with the online learning approach. Inference on 1000 datasets with 2000 posterior samples per parameter took approximately 1.1 seconds.

The results on the SIR model are depicted in Figure 7. In line with the previous examples, we observe very good recovery of the true parameters, with NRMSE at $T = 500$ around 0.03, and R^2 s around 0.99. We observe decent performance even at smaller T s and the expected improvements as T increases. Specifically, the posterior variance shrinks as T increases. The SCB plots indicate that the approximate posteriors are well calibrated, with the approximate posterior mean of β slightly overestimating the true parameter values in the lower range.

3.8 Learned vs. Hand-Crafted Summaries: The Lotka-Volterra Population Model

See Appendix A

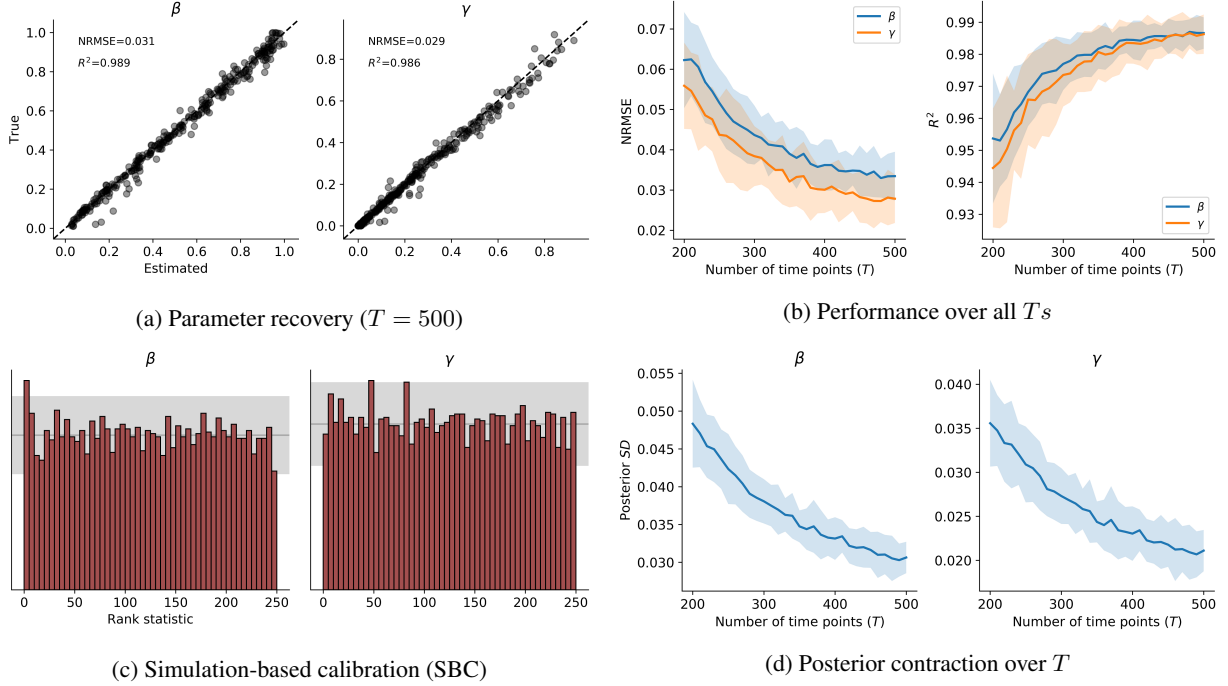


Figure 7: BayesFlow results obtained on the stochastic SIR model.

4 Discussion

In the current work, we proposed and explored a novel method which uses invertible neural networks to perform globally amortized approximate Bayesian inference. The method, which we named BayesFlow, requires only simulations from a forward model to learn an efficient probabilistic mapping between data and parameters. We demonstrated the utility of BayesFlow by applying it to models and data from various research domains. Further, we explored an online learning approach with variable number of observations per iteration. We demonstrated that this approach leads to excellent parameter estimation throughout the examples considered in the current work. In theory, BayesFlow is applicable to any mathematical forward model which can be implemented as a computer simulation. In the following, we highlight the main advantages of BayesFlow.

First, the introduction of separate summary and inference networks renders the method independent of the shape or the size of the observed data. The summary network learns a fixed-size vector representation of the data in an automatic, data-driven manner. Since the summary network is optimized jointly with the inference network, the learned data representation is encouraged to be maximally informative for inferring the parameters' posterior. This is particularly useful in settings where appropriate summary statistics are not known and, as a consequence, relevant information is lost through the choice of sub-optimal summary functions. However, if *sufficient* statistics are available in a given domain, one might omit the summary network altogether and feed these statistics directly to the invertible network.

Second, we showed that BayesFlow generates samples from the correct posterior under perfect convergence without distributional assumptions on the shape of the posterior. This is in contrast to variational methods which optimize a lower-bound on the posterior [26, 24], and oftentimes assume Gaussian approximate posteriors. Additionally, we also showed throughout all examples that the posterior means generated by the BayesFlow method are mostly excellent estimates for the true values. Beyond this, the fact that the BayesFlow method recovers the full posterior over parameters does not necessitate the usage of point estimates or summary statistics of the posterior. Further, we observe the desired posterior contraction (posterior variance decreases with increasing number of observations) and better recovery with increasing number of observations. These are indispensable properties of any Bayesian parameter estimation method, since they mirror the decrease in epistemic uncertainty and the simultaneous increase in information due to availability of more data.

Third, the largest computational cost of BayesFlow is paid during the training phase. Once trained, the networks can efficiently compute the posterior for any observed dataset arising from the forward model. This is similar to the recently introduced *prepaid method* [33]. However, this method memorizes a large database of pre-computed summary

statistics for fast nearest-neighbor inference, whereas a BayesFlow’s network weights define an abstract representation of the relationship between data and parameters over the whole space of hidden parameters. Traditionally, abstract representations like this only existed for analytically invertible model families, whereas more complex forward models required case-based inference, that is, expensive re-training for each observed dataset. Amortized inference as realized by BayesFlow is thus especially advantageous for exploring, testing and comparing competing scientific hypotheses in research domains where an intractable model needs to be fit to multiple independent datasets.

Finally, all computations in the BayesFlow method benefit from a high degree of parallelism and can thus utilize the advantages of modern GPU acceleration.

These advantages notwithstanding, limitations of the proposed method should also be mentioned. Although we could provide a theoretical guarantee that BayesFlow samples from the true joint posterior under perfect convergence, this might not be achieved in practice. Therefore, it is essential that proper calibration of point estimates and estimated joint posteriors is performed for each application of the method. Fortunately, validating a trained BayesFlow architecture is easy due to amortized inference. Below, we discuss potential challenges and limitations of the method.

First, the design of the summary network and inference networks is a crucial choice for achieving optimal performance of the method. As already mentioned, the summary network should be able to represent the observed data without losing essential information and the invertible network should be powerful enough to capture the behavior of the forward model. Nevertheless, in some real-world scenarios, there might be little guidance on how to actually construct suitable summary networks. Recent work on probabilistic symmetry [6] and algorithmic alignment [52] as well as our current experiments do, however, provide some insights about the design of summary networks. For instance, *i.i.d.* data induce a permutation invariant distribution which is well modeled with a deep invariant network [6]. Data with temporal or spatial dependencies are best modeled with recurrent [21], or convolutional [41] networks. When pairwise or multi-way relationships are particularly informative, attention [48] or graph networks [52] appear as reasonable choices. On the other hand, the depth of the invertible network should be tailored to the complexity of the mathematical model of interest. More ACBs will enable the network to encode more complex distributions but will increase training time. Very high-dimensional problems might also require very large networks with millions of parameters, up to a point where estimation becomes practically unfeasible. However, most mathematical models in the life sciences prioritize parsimony and interpretability, so they do not contain hundreds or thousands of latent parameters. In any case, future applications might require novel network architectures and solutions which go beyond our initial recommendations.

Another potential issue is the large number of neural network and optimization hyperparameters that might require fine-tuning by the user for optimal performance on a given task. We observe that excellent performance is often achieved with default settings. Using larger networks consisting of 5 to 10 ACBs does not seem to hurt performance or destabilize training, even if the model to be learned is relatively simple. Based on our results, we expect that a single architecture should be able to perform well on models from a given domain. Future research should investigate this question of generality by applying the method to different or even competing models within different research domains. Future research should investigate the impact of modern hyperparameter optimization methods such as Bayesian optimization [14].

Finally, even though modern deep learning libraries allow for rapid and relatively straightforward development of various neural network architectures, the implementational burden associated with the method is non-trivial. Thus, we are currently developing a general user-friendly software, which will abstract away most intricacies from the users of our method.

We hope that the new BayesFlow method will enable researchers from a variety of fields to accelerate model-based inference and will further prove its utility beyond the examples considered in this paper.

Acknowledgment

We thank Paul Bürkner, Manuel Haussmann, Jeffrey Rouder, Raphael Hartmann, David Izydorczyk, Hannes Wendler, Chris Wendler, and Karin Prillinger for their invaluable comments and suggestions that greatly improved the manuscript. We also thank Francis Tuerlinckx and Stijn Verdonck for their support and thought-provoking ideas.

References

- [1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for large-scale machine learning. In *12th USENIX Symposium on Operating Systems Design and Implementation*, pages 265–283, 2016.

- [2] Lynton Ardizzone, Jakob Kruse, Sebastian Wirkert, Daniel Rahner, Eric W Pellegrini, Ralf S Klessen, Lena Maier-Hein, Carsten Rother, and Ullrich Köthe. Analyzing inverse problems with invertible neural networks. In *Intl. Conf. on Learning Representations*, 2019.
- [3] Lynton Ardizzone, Carsten Lüth, Jakob Kruse, Carsten Rother, and Ullrich Köthe. Guided image generation with conditional invertible neural networks. *arXiv:1907.02392*, 2019.
- [4] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019.
- [5] Daniele Bigoni, Olivier Zahm, Alessio Spantini, and Youssef Marzouk. Greedy inference with layers of lazy maps. *arXiv:1906.00031*, 2019.
- [6] Benjamin Bloem-Reddy and Yee Whye Teh. Probabilistic symmetry and invariant neural networks. *arXiv:1901.06082*, 2019.
- [7] William M Bolstad and James M Curran. *Introduction to Bayesian statistics*. John Wiley & Sons, 2016.
- [8] Laming Chen, Guoxin Zhang, and Eric Zhou. Fast greedy map inference for determinantal point process to improve recommendation diversity. In *Advances in Neural Information Processing Systems*, pages 5622–5633, 2018.
- [9] Kyle Cranmer, Johann Brehmer, and Gilles Louppe. The frontier of simulation-based inference. *arXiv:1911.01429*, 2019.
- [10] Katalin Csilléry, Michael GB Blum, Oscar E Gaggiotti, and Olivier François. Approximate Bayesian Computation (ABC) in Practice. *Trends in Ecology & Evolution*, 25(7):410–418, 2010.
- [11] Matthew C Deans. Maximally informative statistics for localization and mapping. In *Proceedings 2002 IEEE International Conference on Robotics and Automation (Cat. No. 02CH37292)*, volume 2, pages 1824–1829. IEEE, 2002.
- [12] Gianluca Detommaso, Jakob Kruse, Lynton Ardizzone, Carsten Rother, Ullrich Köthe, and Robert Scheichl. HINT: hierarchical invertible neural transport for general and sequential bayesian inference. *arXiv:1905.10687*, 2019.
- [13] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using Real NVP. *arXiv:1605.08803*, 2016.
- [14] Katharina Eggenberger, Matthias Feurer, Frank Hutter, James Bergstra, Jasper Snoek, Holger Hoos, and Kevin Leyton-Brown. Towards an empirical foundation for assessing bayesian optimization of hyperparameters. In *NIPS workshop on Bayesian Optimization in Theory and Practice*, volume 10, page 3, 2013.
- [15] Felix A Gers, Jürgen Schmidhuber, and Fred Cummins. Learning to forget: Continual prediction with LSTM. *Neural Computation*, 12(19):2451–2471, 2000.
- [16] David S Greenberg, Marcel Nonnenmacher, and Jakob H Macke. Automatic posterior transformation for likelihood-free inference. *arXiv:1905.07488*, 2019.
- [17] Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(Mar):723–773, 2012.
- [18] Aditya Grover, Manik Dhar, and Stefano Ermon. Flow-gan: Combining maximum likelihood and adversarial learning in generative models. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [19] John R Hershey and Peder A Olsen. Approximating the kullback leibler divergence between gaussian mixture models. In *2007 IEEE International Conference on Acoustics, Speech and Signal Processing-ICASSP’07*, volume 4, pages IV–317. IEEE, 2007.
- [20] Herbert W Hethcote. The mathematics of infectious diseases. *SIAM review*, 42(4):599–653, 2000.
- [21] Seong Jae Hwang, Zirui Tao, Won Hwa Kim, and Vikas Singh. Conditional recurrent flow: Conditional generation of longitudinal samples with applications to neuroimaging. *arXiv:1811.09897*, 2018.
- [22] Bai Jiang, Tung-yu Wu, Charles Zheng, and Wing H Wong. Learning summary statistic for approximate bayesian computation via deep neural network. *Statistica Sinica*, pages 1595–1618, 2017.
- [23] Matt J Keeling and Pejman Rohani. *Modeling infectious diseases in humans and animals*. Princeton University Press, 2011.
- [24] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *stat*, 1050:1, 2014.

- [25] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In *Advances in Neural Information Processing Systems*, pages 10215–10224, 2018.
- [26] Durk P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. Improved variational inference with inverse autoregressive flow. In *Advances in Neural Information Processing Systems*, pages 4743–4751, 2016.
- [27] Emmanuel Klinger, Dennis Rickert, and Jan Hasenauer. pyabc: distributed, likelihood-free inference. *Bioinformatics*, 34(20):3591–3593, 2018.
- [28] Hongzhou Lin and Stefanie Jegelka. Resnet with one-neuron hidden layers is a universal approximator. In *Advances in Neural Information Processing Systems*, pages 6169–6178, 2018.
- [29] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pages 3431–3440, 2015.
- [30] Jan-Matthis Lueckmann, Pedro J Goncalves, Giacomo Bassetto, Kaan Öcal, Marcel Nonnenmacher, and Jakob H Macke. Flexible statistical inference for mechanistic models of neural dynamics. In *Advances in Neural Information Processing Systems*, pages 1289–1299, 2017.
- [31] Ulf Kai Mertens. *Deep learning methods for likelihood-free inference: approximating the posterior distribution with convolutional neural networks*. PhD thesis, Heidelberg University, 2019.
- [32] Ulf Kai Mertens, Andreas Voss, and Stefan Radev. Abrox—a user-friendly python module for approximate bayesian computation with a focus on model comparison. *PloS one*, 13(3):e0193981, 2018.
- [33] Merijn Mestdag, Stijn Verdonck, Kristof Meers, Tim Loossens, and Francis Tuerlinckx. Prepaid parameter estimation without likelihoods. *PLoS computational biology*, 15(9):e1007181, 2019.
- [34] Steven Miletić, Brandon M Turner, Birte U Forstmann, and Leendert van Maanen. Parameter recovery for the leaky competing accumulator model. *Journal of Mathematical Psychology*, 76:25–50, 2017.
- [35] Ashish Mishra, Shiva Krishna Reddy, Anurag Mittal, and Hema A Murthy. A generative model for zero shot learning using conditional variational autoencoders. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 2188–2196, 2018.
- [36] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.
- [37] George Papamakarios and Iain Murray. Fast ϵ -free inference of simulation models with bayesian conditional density estimation. In *Advances in Neural Information Processing Systems*, pages 1028–1036, 2016.
- [38] George Papamakarios, David C Sterratt, and Iain Murray. Sequential neural likelihood: Fast likelihood-free inference with autoregressive flows. *arXiv:1805.07226*, 2018.
- [39] Mijung Park, Wittawat Jitkrittum, and Dino Sejdinovic. K2-ABC: approximate bayesian computation with kernel embeddings. In *Intl. Conf. Artificial Intelligence and Statistics*, pages 398–407, 2016.
- [40] Matthew D Parno and Youssef M Marzouk. Transport map accelerated markov chain monte carlo. *SIAM/ASA Journal on Uncertainty Quantification*, 6(2):645–682, 2018.
- [41] Stefan T Radev, Ulf K Mertens, Andreas Voss, and Ullrich Köthe. Towards end-to-end likelihood-free inference with convolutional neural networks. *British Journal of Mathematical and Statistical Psychology*, 2019.
- [42] Roger Ratcliff and Gail McKoon. The diffusion decision model: theory and data for two-choice decision tasks. *Neural computation*, 20(4):873–922, 2008.
- [43] Louis Raynal, Jean-Michel Marin, Pierre Pudlo, Mathieu Ribatet, Christian P Robert, and Arnaud Estoup. ABC random forests for Bayesian parameter inference. *Bioinformatics*, 35(10):1720–1728, 2018.
- [44] Scott A Sisson and Yanan Fan. *Likelihood-free MCMC*. Chapman & Hall/CRC, New York.[839], 2011.
- [45] Mikael Sunnåker, Alberto Giovanni Busetto, Elina Numminen, Jukka Corander, Matthieu Foll, and Christophe Dessimoz. Approximate bayesian computation. *PLoS computational biology*, 9(1):e1002803, 2013.
- [46] Sean Talts, Michael Betancourt, Daniel Simpson, Aki Vehtari, and Andrew Gelman. Validating bayesian inference algorithms with simulation-based calibration. *arXiv:1804.06788*, 2018.
- [47] Brandon M Turner and Per B Sederberg. A generalized, likelihood-free method for posterior estimation. *Psychonomic bulletin & review*, 21(2):227–250, 2014.
- [48] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.

- [49] Andreas Voss, Veronika Lerche, Ulf Mertens, and Jochen Voss. Sequential sampling models with variable boundaries and non-normal noise: A comparison of six models. *Psychonomic bulletin & review*, pages 1–20, 2019.
- [50] Darren J Wilkinson. *Stochastic modelling for systems biology*. CRC press, 2011.
- [51] Simon N Wood. Statistical inference for noisy nonlinear ecological dynamic systems. *Nature*, 466(7310):1102, 2010.
- [52] Keyulu Xu, Jingling Li, Mozhi Zhang, Simon S Du, Ken-ichi Kawarabayashi, and Stefanie Jegelka. What can neural networks reason about? *arXiv:1905.13211*, 2019.

Appendix

A Learned vs. Hand-Crafted Summaries: The Lotka-Volterra Population Model

With this final example, we want to compare the performance of our method with an LSTM summary network vs. performance obtained with a standard set of hand-crafted summary statistics. For this, we focus on the well-studied Lotka-Volterra (LV) model. The LV model describes the dynamics of biological systems in which a population of predators interacts with a population of prey [50]. It involves a pair of first order, non-linear, differential equations given by:

$$\frac{d}{dt} = \alpha u - \beta uv \quad (1)$$

$$\frac{d}{dt} = -\gamma v + \delta \beta uv \quad (2)$$

where u denotes the number of preys, v denotes the number of predators, and the parameter vector controlling the interaction between the species is $\theta = (\alpha, \beta, \gamma, \delta)$.

During training of the networks, we set the initial conditions as $u_0 = 10$ and $v_0 = 5$ and consider an interval $I_T = 15$ of discrete time units with $T = 500$ time steps (samples) in between. Each sample x_t in each LV time-series $x_{1:T}$ is thus a 2-dimensional vector containing the number of prey and predators in the population at time unit t .

We train two invertible neural networks. The first is trained jointly with an LSTM summary network which outputs a 9-dimensional learned summary statistic $h_\psi(x_{1:T})$. The second uses a set of 9 typically used, hand-crafted summary statistics [37, 38], which include: the mean of the time series; the log variance of the time-series; the auto-correlation of each timeseries at lags 0.2 and 0.4 time units; the cross-correlation between the two time series. The same cINN architecture with 5 ACBs is used for both training scenarios. For each scenario, we perform the same number of iterations and epochs. Online learning for each training scenario took approximately 4 hours in total wall-clock time.

The results obtained on the LV model are depicted in Figure S1. We observe notably better recovery of the true parameter estimates when performing inference with the learned summary statistics. The approximate posteriors are also better calibrated when conditioned on the set of 9 learned summary statistics. These results highlight the advantages of using a summary networks when no sufficient summary statistics are available. Finally, Figure S1e and Figure S1f depict the posteriors obtained by the two different INNs on a single dataset with ground-truth parameters $\theta = (1, 1, 1, 1)$. Evidently, learning the summary statistics leads to much sharper posteriors and better point-estimate recovery.

B Computation of Validation Metrics

Normalized Root Mean Squared Error

The normalized root mean squared error (NRMSE) between a sample of true parameters $\{\theta^{(m)}\}_{m=1}^M$ and a sample of estimated parameters $\{\hat{\theta}^{(m)}\}_{m=1}^M$ is given by:

$$NRMSE = \sqrt{\sum_{m=1}^M \frac{(\theta^{(m)} - \hat{\theta}^{(m)})^2}{\theta_{max} - \theta_{min}}} \quad (3)$$

Due to the normalization factor $\theta_{max} - \theta_{min}$, the NRMSE is scale-independent, and thus suitable for comparing the recovery across parameters with different numerical ranges. The NRMSE equals zero when the estimates are exactly equal to the true values.

Coefficient of Determination

The coefficient of determination R^2 measures the proportion of variance in a sample of true parameters $\{\theta^{(m)}\}_{m=1}^M$ that is *explained* by a sample of estimated parameters $\{\hat{\theta}^{(m)}\}_{m=1}^M$. It is computed as:

$$R^2 = 1 - \sum_{m=1}^M \frac{(\theta^{(m)} - \hat{\theta}^{(m)})^2}{(\theta^{(m)} - \bar{\theta}^{(m)})^2} \quad (4)$$

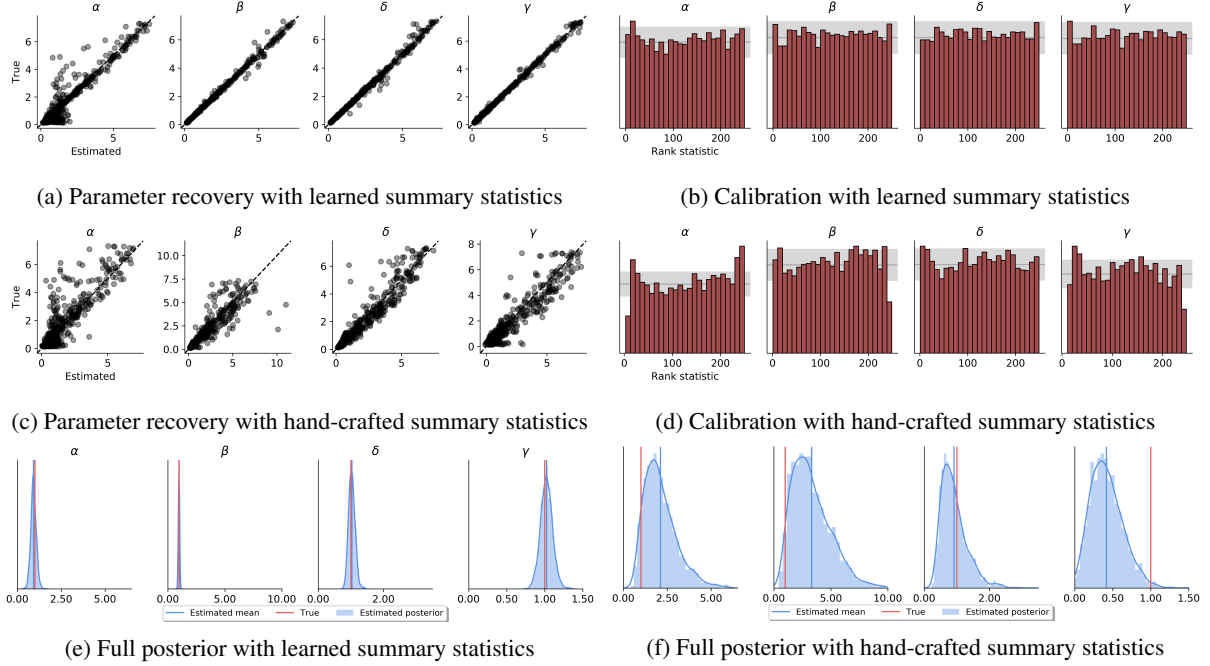


Figure S1: Comparison of recovery/calibration on the LV model with learned vs. hand-crafted summary statistics **(a)** Simulation-based calibration (SBC) with learned summary statistics; **(b)** Parameter recovery with learned summary statistics; **(c)** Parameter recovery with hand-crafted summary statistics; **(d)** Simulation-based calibration (SBC) with hand-crafted summary statistics; **(e)** Example full posteriors obtained on a single dataset with ground-truth parameters $\theta = (1, 1, 1, 1)$ obtained with learned summaries; **(f)** The posterior obtained from the same dataset using hand-crafted summary statistics.

where $\bar{\theta}$ denotes the mean of the true parameter samples. When R^2 equals 1, the estimates are perfect reconstructions of the true parameters.

Re-simulation Error

To compute the re-simulation error Err_{sim} , we first obtain an estimate of the true parameter value given an observed (validations) dataset $\mathbf{x}_{1:N}^o$ by computing the mean of the approximate posterior $\tilde{\theta}$. Then, we run the mathematical model to obtain a simulated dataset $\mathbf{x}_{1:N}^s = g(\tilde{\theta}, \xi)$. Finally, we compute the maximum mean discrepancy (MMD, [17]) between the observed and the simulated dataset $MMD(\mathbf{x}_{1:N}^o, \mathbf{x}_{1:N}^s)$. The MMD is a kernel-based metric which estimates the mismatch between two distributions given samples from the distributions by comparing all of their moments. It equals zero when the two distributions are equal almost everywhere [17]). Thus, a low MMD indicates that the distribution of $\mathbf{x}_{1:N}^s$ is close to the distribution of $\mathbf{x}_{1:N}^o$. Conversely, a high MMD indicates that the distribution of $\mathbf{x}_{1:N}^s$ is far from the distribution of $\mathbf{x}_{1:N}^o$. We report the median MMD computed over all validation datasets.

Calibration Error

The calibration error Err_{cal} quantifies how well the coverage of an approximate posterior matches the coverage of an unknown true posterior. Let α_θ be the fraction of true parameter values lying in a corresponding α -credible interval of the approximate posterior. Thus, for a perfectly calibrated approximate posterior, α_θ should equal α for all $\alpha \in (0, 1)$. We compute the calibration error for each marginal posterior as the median absolute deviation $|\alpha_\theta - \alpha|$ for 100 equally spaced values of $\alpha \in (0, 1)$. Therefore, the calibration error ranges between 0 and 1 with 0 indicating perfect calibration and 1 indicating complete miscalibration of the approximate posterior.

Kullback-Leibler Divergence

The Kullback-Leibler divergence (\mathbb{KL}) quantifies the increase in entropy incurred by approximating a target probability distribution P with a distribution Q . Its general form for absolutely continuous distributions is given by

$$\mathbb{KL}(P \parallel Q) = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx \quad (5)$$

where p and q denote the pdfs of P and Q . In the case where P and Q are both multivariate Gaussian distributions, the KL divergence can be computed in closed form [19]:

$$\mathbb{KL}(P \parallel Q) = \frac{1}{2} \left[\log \frac{\det \Sigma_q}{\det \Sigma_p} + \text{Tr}(\Sigma_q^{-1} \Sigma_p) - d + (\mu_p - \mu_q)^T \Sigma_q^{-1} (\mu_p - \mu_q) \right] \quad (6)$$

where Σ_p and Σ_q denote the covariance matrices of p and q , μ_p and μ_q the respective mean vectors, and d the number of dimensions of the Gaussian. In the case of diagonal Gaussian distributions, Eq.6 reduces to:

$$\mathbb{KL}(P \parallel Q) = \sum_{i=1}^d \left(\log \frac{\sigma_{q,i}}{\sigma_{p,i}} + \frac{\sigma_{p,i}^2 + (\mu_{q,i} - \mu_{p,i})^2}{2\sigma_{q,i}^2} - \frac{1}{2} \right) \quad (7)$$

Even though the KL divergence is not a proper distance metric, as it is not symmetric in its arguments, it can be used to quantify the error of approximation when a closed-form solution is available.

Simulation-Based Calibration

Simulation-based calibration is a method to detect systematic biases in any Bayesian posterior sampling method [46]. It is based on the *self-consistency* of the Bayesian joint distribution. Given a sample from the prior distribution $\tilde{\theta} \sim p(\theta)$ and a sample from the forward model $\tilde{x} \sim p(x \mid \tilde{\theta})$, one can integrate $\tilde{\theta}$ and \tilde{x} out of the joint distribution and recover back the prior of θ :

$$p(\theta) = \int p(\theta, \tilde{\theta}, \tilde{x}) d\tilde{x} d\tilde{\theta} \quad (8)$$

$$= \int p(\theta, \tilde{x} \mid \tilde{\theta}) p(\tilde{\theta}) d\tilde{x} d\tilde{\theta} \quad (9)$$

$$= \int p(\theta \mid \tilde{x}) p(\tilde{x} \mid \tilde{\theta}) p(\tilde{\theta}) d\tilde{x} d\tilde{\theta} \quad (10)$$

If the Bayesian sampling method produces samples from the exact posterior, the equality implied by Eq.10 should hold regardless of the particular form of the posterior. Thus, any violation of this equality indicates some error incurred by the sampling method. The authors of [46] propose **Algorithm 2** for visually detecting such violations:

Algorithm 2 Simulation-based calibration (SBC) for a single parameter θ

- 1: **for** $m = 1, \dots, M$ **do**
 - 2: Sample $\tilde{\theta}^{(m)} \sim p(\theta)$
 - 3: Simulate a dataset $\mathbf{x}_{1:N}^{(m)} = g(\tilde{\theta}^{(m)}, \boldsymbol{\xi})$
 - 4: Draw posterior samples $\{\theta^{(l)}\}_{l=1}^L \sim p_\phi(\theta \mid \mathbf{x}_{1:N}^{(m)})$
 - 5: Compute rank statistic $r^{(m)} = \sum_{l=1}^L \mathbb{1}_{[\theta^{(l)} < \tilde{\theta}^{(m)}]}$
 - 6: Store $r^{(m)}$
 - 7: **end for**
 - 8: Create a histogram of $\{r^{(i)}\}_{i=1}^M$ and inspect it for uniformity
-

Algorithm 2 is correct, since Eq.10 implies that the rank statistic defined in line 5 should be uniformly distributed. Hence, any deviations from uniformity indicate some interpretable error in the approximate posterior [46].

C Model Details

The Ricker Model

Summary Network. We use a bidirectional long short-term memory (LSTM) recurrent neural network [15] for the raw Ricker time-series. The LSTM network architecture is a reasonable choice for this example, as it is able to capture long-term dependencies in datasets with temporal or spatial autocorrelations. LSTMs can also easily deal with variable-length time-series.

Simulation. We place the following uniform priors over the Ricker model parameters:

$$\rho \sim \mathcal{U}(0, 15) \quad (11)$$

$$r \sim \mathcal{U}(1, 90) \quad (12)$$

$$\sigma \sim \mathcal{U}(0.05, 0.7) \quad (13)$$

These ranges appear to be very broad, as datasets generated by extreme parameter values appear implausible in real-world scenarios. Nevertheless, we stick to broad priors for training, even though parameter recovery might degrade at the extremes.

Figure S2 depicts different simulated Ricker timeseries generated via draws from the prior.

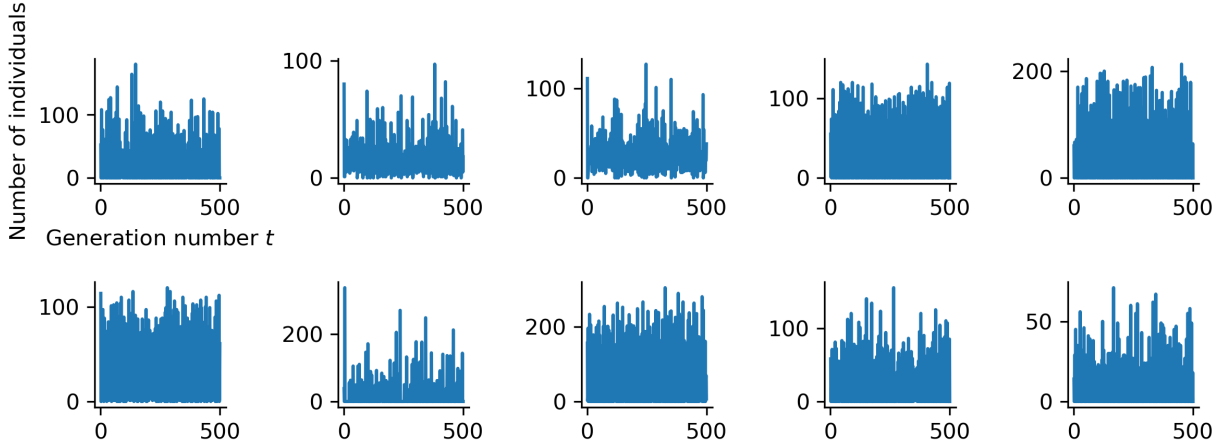


Figure S2: Example Ricker datasets generated with different parameters.

The Lévy-Flight Model

Summary Network. We use a permutation invariant neural network [6] for the *i.i.d.* reaction times (RT) data. Similarly to the toy Regression example, each response in an RT dataset is assumed to be independent of all others, so permutations of the dataset must lead to the same parameter estimates.

Simulation. We place the following uniform priors over the LFM parameters, since they are broad enough to cover the range of realistic RT distributions encountered in empirical choice RT scenarios:

$$v_0 \sim \mathcal{U}(0, 6) \quad (14)$$

$$v_1 \sim \mathcal{U}(-6, 0) \quad (15)$$

$$zr \sim \mathcal{U}(0.3, 0.7) \quad (16)$$

$$a \sim \mathcal{U}(0.6, 3) \quad (17)$$

$$t_0 \sim \mathcal{U}(0.3, 1) \quad (18)$$

$$\alpha \sim \mathcal{U}(1, 2) \quad (19)$$

Figure S3 depicts different simulated RT distributions generated via draws from the prior.

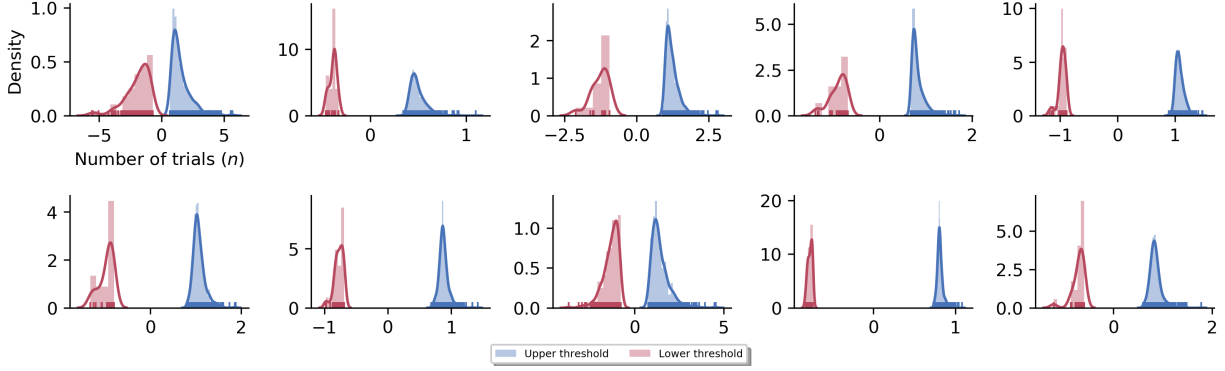


Figure S3: Example RT distributions generated with different parameters.

The Stochastic SIR Model

Summary Network. We use a 1D fully convolutional neural network [29] for the raw SIR time-series into fixed-size vectors. Here, we choose a convolutional network architecture over the previously mentioned LSTM, as convolutional networks are more computationally efficient. Further, we wanted to underline the utility of 1D convolutional networks for multidimensional time-series data. Finally, convolutional networks can also deal with variable input sizes.

Simulation. We place the following uniform priors over the two rate parameters of the stochastic SIR model:

$$\beta \sim \mathcal{U}(0.01, 1) \quad (20)$$

$$\gamma \sim \mathcal{U}(0.01, \beta) \quad (21)$$

These ranges were chosen based on empirical plausibility of the generated SIR time-series.

Figure S4 depicts different SIR timeseries generated via draws from the prior.

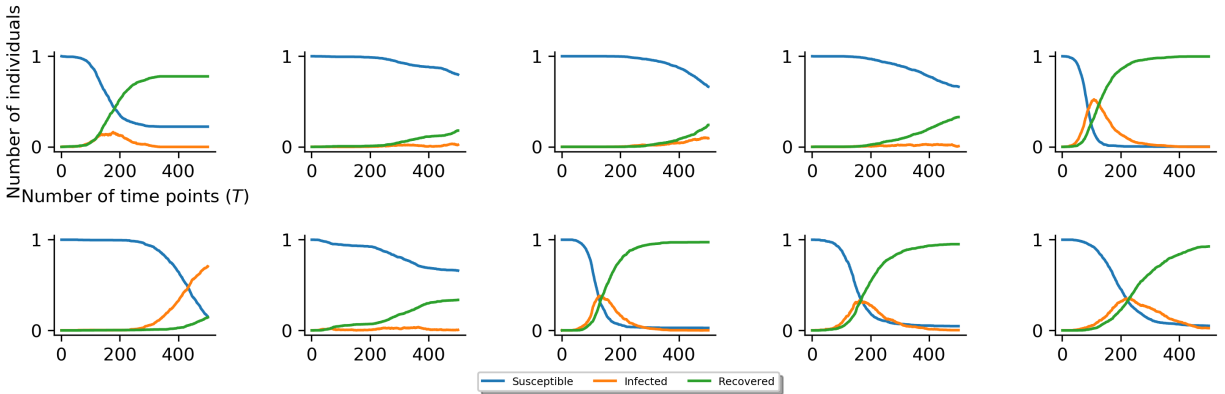


Figure S4: Example SIR timeseries generated with different parameters.

The Lotka-Volterra Model

Summary Network. We use a bidirectional long short-term memory (LSTM) recurrent neural network [15] for the raw LV time-series (as in the Ricker example).

Simulation. We place the following broad uniform priors over the LV parameters. Some of the parameter combinations produced divergent simulations, which we removed during online learning.

$$\alpha \sim \mathcal{U}(\exp(-2), \exp(2)) \quad (22)$$

$$\beta \sim \mathcal{U}(\exp(-2), \exp(2)) \quad (23)$$

$$\gamma \sim \mathcal{U}(\exp(-2), \exp(2)) \quad (24)$$

$$\delta \sim \mathcal{U}(\exp(-2), \exp(2)) \quad (25)$$

D Example Posteriors on Ricker Datasets

Marginal posteriors from ten validation datasets simulated from the Ricker model are depicted in Figure S5. We observe widely different posterior shapes, highlighting the importance of working with arbitrary posterior shapes.

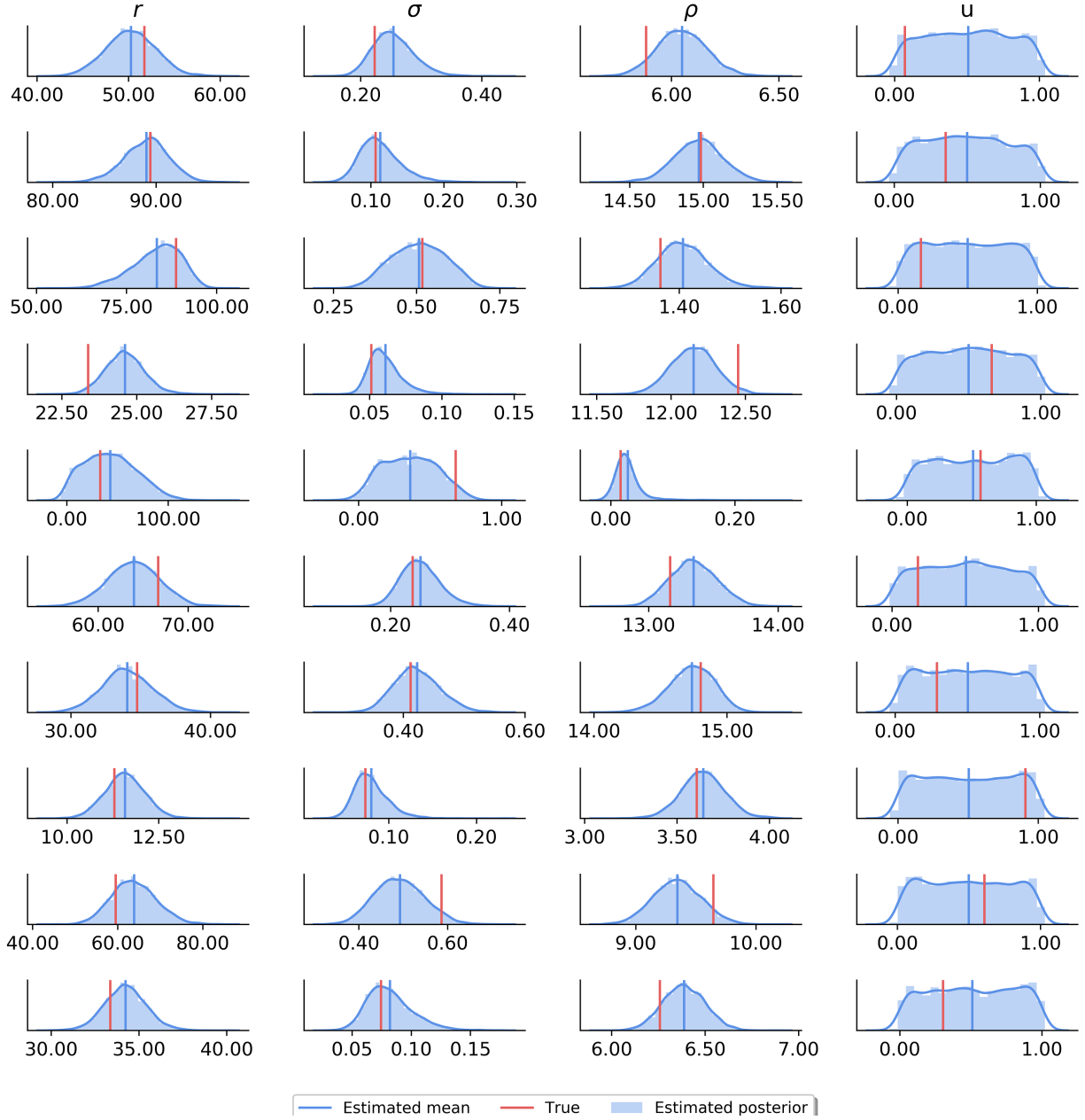


Figure S5: Ten example Ricker marginal posteriors