

# Learning complex stochastic models with invertible neural networks: a likelihood-free Bayesian approach

Stefan T. Radev<sup>1</sup>, Ulf K. Mertens<sup>1</sup>, Andreas Voss<sup>1</sup>, Lynton Ardizzone<sup>2</sup>, and Ullrich Köthe<sup>2</sup>

<sup>1</sup> Institute of Psychology, Heidelberg University, Hauptstr. 47-51, 69117 Heidelberg, Germany; <sup>2</sup> Heidelberg Collaboratory for Image Processing (HCI), Interdisciplinary Center for Scientific Computing (IWR), Heidelberg University, Im Neuenheimer Feld 205, 69120 Heidelberg, Germany

This manuscript was compiled on July 11, 2019

1 Mathematical models of complex processes are ubiquitous throughout the sciences. As the processes under study and the models de-  
2 scribing them become increasingly complex, parameter estimation with standard Bayesian and frequentist methods can quickly become  
3 intractable. To address this issue, we propose a novel method for likelihood-free inference based on invertible neural networks. The method  
4 is capable of performing full Bayesian inference on large datasets by training the networks on simulated data to learn a probabilistic map-  
5 ping between parameters and data. The method is independent of any particular data representation, as it incorporates a summary network  
6 trained to embed the observed data into fixed-size vectors in a data-driven way. This makes the method applicable to various scenarios where  
7 standard inference techniques fail. We demonstrate the utility of the method on a toy model with known analytic posterior and on example  
8 models from population dynamics, epidemiology, cognitive science and genetics. We argue that our method provides a general framework  
9 for building reusable parameter estimation machines for any process model from which data can be simulated.

Deep learning | Invertible networks | Bayesian inference | Parameter estimation | Stochastic models

**M**athematical models are formal descriptions of scientific theories. In its most abstract form, a mathematical model is specified by a set of parameters  $\theta$  and a forward model  $q$  mimicking the process by which manifest data  $x$  arise from latent parameters:

$$x = q(\theta) \quad [1]$$

1 The functional form of  $q$ , which can represent an arbitrarily complex process by an arbitrarily complicated expression, is usually guided by a well-founded theoretical framework. For instance,  $q$  can be a stochastic differential equation describing the  
2 dynamics of single neurons in the brain, or a step-by-step biological algorithm controlling the rate of gene expression in certain  
3 cells. Thus, it is only within the context of a theory that a meaningful interpretation in terms of some mechanism can be attached to the parameters of a mathematical model. Examples of mathematical models can be found in various scientific  
4 domains, for instance, genetics (1, 2), cognitive science (3, 4), neuroscience (5, 6), population dynamics (7, 8), epidemiology  
5 (9, 10), to name just a few.

6 An important goal of mathematical modeling is recovering the parameters of the model from observed data. Idealized  
7 parameter estimation involves computing the inverse (backward) model  $\theta = q^{-1}(x)$  exactly. Unfortunately, this inverse rarely exists,  
8 and estimating parameters remains computationally inconvenient in a wide range of settings.

9 Parameter estimation becomes notably difficult when the forward model in Eq. 1 does not provide a closed-form for the  
10 likelihood function (3, 11, 12). This poses great difficulties for Bayesian and frequentist methods alike, since both depend on  
11 the numerical evaluation of a likelihood function as a proxy for assessing model fit to data. Even if a likelihood function is available in closed-form, inference may be prohibitively slow for real-world applications. In this case, enforcing simplifying  
12 distributional assumptions (i.e., independence or Gaussian assumptions) can increase computational speed, but can also lead to  
13 model misspecifications and dramatically incorrect estimates. Therefore, there is a need for accurate and computationally  
14 convenient likelihood-free estimation methods.

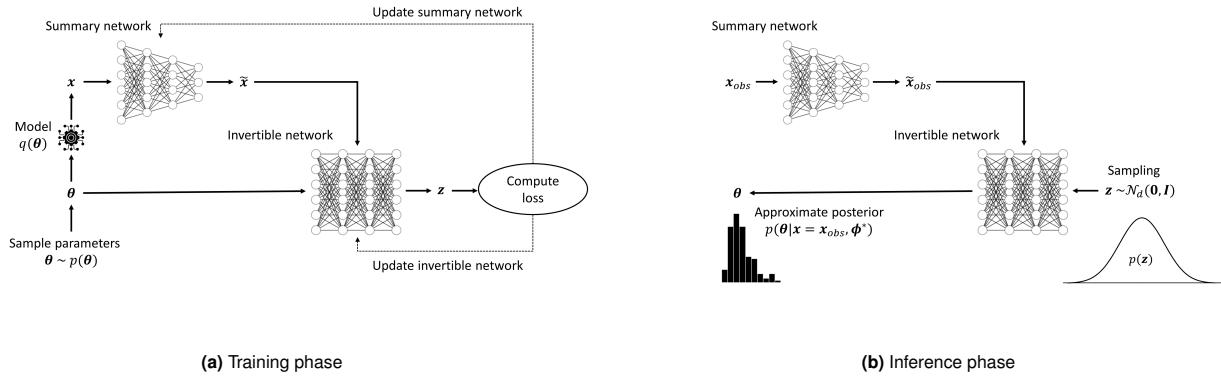
## Significance Statement

Describing complex stochastic processes with parametric models lies at the heart of science. Simulating models given a set of parameters is relatively easy with the aid of modern computers, but inferring model parameters from observed data can often be a challenging endeavor. We combine recent advances in deep learning and Bayesian inference into a powerful method for building reusable parameter estimation networks applicable to various types of models and data encountered in different research fields.

Please provide details of author contributions here.

This research was supported by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation; grant number GRK 2277 "Statistical Modeling in Psychology")

<sup>2</sup>To whom correspondence should be addressed. E-mail: stefan.raeve@psychologie.uni-heidelberg.de



**Fig. 1.** Graphical illustration of the method. **(a)** During the training phase, the summary and the invertible network are trained on simulated data from the model and updated after each batch of simulations; **(b)** During the inference phase, the true posterior of the model parameters is approximated from real data using the trained networks. Thus, knowledge about the relationship between parameters and data (the mathematical model) is compactly encoded within the weights of the two networks. The trained networks can then be shared and used across researchers working on the same model.

Likelihood-free methods bypass the above problems by resorting to a simulation-based approach to parameter estimation and model selection (3, 13). A subset of likelihood-free methods includes approximate Bayesian computation (ABC) methods, which aim at preserving the advantages of Bayesian data analysis even when the likelihood function is intractable or practically impossible to compute (11, 13, 14). ABC methods approximate the likelihood function by repeatedly sampling parameters from a pre-specified prior distribution  $p(\theta)$  and then simulating multiple datasets by running the forward model  $q(\theta)$  using the sampled parameters. Correspondingly, the core ingredients of ABC methods are a prior on  $\theta$ , and a generative model  $q(\theta)$ , usually specified as a function code in a general-purpose programming language (11, 15).

Performing approximate inference comes at the cost of incurring additional approximation error, which accumulates on top of the irreducible estimation error. Within the context of approximate inference, the most common manifestations of approximation error include: *i*) imprecise form of the posterior; *ii*) imprecise posterior moments; *iii*) under- or overestimation of uncertainty. Different approximation methods usually involve multiple trade-offs between minimizing approximation error and keeping computational time within reasonable bounds (3, 16).

To address the shortcomings of traditional methods, ideas from machine learning and deep learning research have recently entered the field of likelihood-free inference (5, 6, 17–22). 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  $D = \{h(x^{(i)}), \theta^{(i)}\}_{i=1}^n$  is created by repeatedly sampling from  $p(\theta)$  and simulating an artificial dataset  $x$  by running  $q(\theta)$  with the sampled parameters. Usually, the dimensionality of the simulated data is reduced by computing summary statistics with a fixed summary function  $h(x)$ . Then, a supervised learning algorithm  $f(h(x); \phi) = \hat{\theta}$  with learnable parameters  $\phi$  (e.g., linear regression, random forest, neural network) is trained on the simulated data to output an estimate of the true data generating parameters. Thus,  $f(h(x); \phi)$  essentially attempts to “learn” the intractable inverse model  $\theta = q^{-1}(x)$ .

Inspired by previous machine learning approaches (5, 6, 19–22), the current work proposes a novel and universal likelihood-free method capable of performing full Bayesian inference on any mathematical process model from which simulations can be obtained. It treats parameter inference as a task of inverting the forward model in Eq. 1 and achieves this by drawing on the modern framework of deep probabilistic modeling for tackling intractable posteriors (23–26). The method integrates two separate deep neural networks modules (detailed in the **Methods** section; see Figure 1) trained jointly on simulated data: a *summary network* and an *invertible network*.

The *summary network* is responsible for learning the most informative summary statistics directly from data. It should be designed to follow the functional and probabilistic symmetries inherent in the data, e.g. a permutationally invariant network for *i.i.d.* data (27), a recurrent network for time-series data (28), or a convolutional network for grid-like data (29). The computation of summary statistics is a crucial aspect in likelihood-free inference. Previous approaches typically used hand-crafted summary statistics tailored to the specific application (14, 20, 21). However, in many applications, finding an optimal set of summary statistics is far from obvious. The summary network eliminates the need to manually decide on and compute a fixed set of summary statistics and, thus, makes the method independent of the format or the size of the data.

The *invertible network* is responsible for learning the posterior of the model parameters given the observed and summarized data. It is based on the recently developed flow-based architecture (24–26). Flow-based methods provide exact latent-variable inference and log-likelihood evaluation when operating at optimum. In the **Methods** section, we show that our method maximizes the posterior over model parameters directly when cast in the context of likelihood-free inference. Furthermore, flow-based methods are capable of approximating high-dimensional distributions (e.g., the pixels of an image). Once trained with a sufficient amount of simulated data, the invertible network can perform rapid Bayesian inference on large datasets from a given research domain.

The joint training of a summary network and an invertible network results in a powerful and universal parameter estimation

59 machine capable of estimating complex mathematical models in various scientific domains (Figure 1). Moreover, the method  
 60 addresses many of the limitations of previous likelihood-free methods. First, it involves no costly MCMC or rejection sampling,  
 61 which makes inference lightning fast, once the networks have been trained. Second, it involves no fixed summary statistics, but  
 62 instead learns the most informative representation of the data in an end-to-end manner. Third, the method is fully Bayesian,  
 63 as it directly learns the posterior over model parameters and thus allows for the quantification of uncertainty, which is a crucial  
 64 requirement in parameter estimation (30, 31). Last, the trained networks can be shared and reused by multiple researchers  
 65 within a scientific domain, thus removing the need for wasteful computations and fitting a separate model for each and every  
 66 dataset. This pooling of computational resources across researchers is an important step forward in mathematical modeling  
 67 (19).

68 To illustrate the utility of the new method, we first apply it to a toy Bayesian regression model with known posterior. Then,  
 69 we present applications to intractable models from cognitive science, population dynamics, epidemiology, and genetics and  
 70 demonstrate state-of-the art parameter recovery. Across the examples, we introduce multiple tools to validate the performance  
 71 of our method. The outline of the remaining manuscript is as follows: The **Methods** section introduces the main building  
 72 blocks of the new method and summarizes the main steps in pseudocode. The **Results** section presents the various applications  
 73 of the model to real-world research domains. Finally, the **Discussion** section lists the advantages of the current method, treats  
 74 some potential pitfalls and explores future research vistas. Python code and simulation scripts for all current applications are  
 75 freely available as Jupyter notebooks at <https://github.com/stefanradev93/cINN> and as a small library based on *TensorFlow*  
 76 (32) for creating and training custom invertible networks with GPU support, along with some validation tools.

## 77 Methods

78 **Notation.** In the following, we denote observed or simulated univariate datasets from the mathematical model of interest as  
 79  $\mathbf{x} = (x_1, x_2, \dots, x_n)$ , and multivariate datasets as  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$ . The parameters of a mathematical model are represented  
 80 as a vector  $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_d)$ , and all trainable parameters of the invertible and summary neural networks as  $\phi = (\phi_{inv}, \phi_{sum})$ .  
 81 The number of parameters of a mathematical model will be denoted as  $d$ , and the number of observations in  $\mathbf{x}$  or  $\mathbf{X}$  as  $n$ .

82 **Deep Probabilistic Modeling.** Our method draws on major advances in modern deep probabilistic modeling, also referred to  
 83 as deep generative modeling (23, 24, 27, 33). A hallmark idea in deep probabilistic modeling is to handle intractable target  
 84 probability distributions by sampling from simpler distributions (e.g., Gaussian or uniform distributions) and transforming these  
 85 samples via learned complex non-linear transformations. Most popular deep probabilistic models entail two phases. During the  
 86 *training phase*, a transformation from the simple to the desired target distribution is learned by optimizing a cost function via  
 87 backpropagation (see Figure 1a). During the *inference phase*, samples from the target distribution are obtained by sampling  
 88 from the simple distribution and applying the transformation learned during the training phase (see Figure 1b). Using this  
 89 approach, recent applications of deep probabilistic models have achieved unprecedented results on extremely high-dimensional  
 90 and intractable problems (e.g., complex data distributions such as images, music, or text) (24, 25, 27).

91 In the context of mathematical modeling and Bayesian inference, the target distribution is the posterior distribution of  
 92 model parameters  $p(\boldsymbol{\theta}|\mathbf{x})$  capturing the uncertainty about the numerical values of parameters. We can leverage the fact that  
 93 most mathematical models are generative in nature and as such can be used to perform multiple simulations of the process of  
 94 interest. By specifying a prior distribution over the model parameters  $p(\boldsymbol{\theta})$ , one can generate arbitrarily large datasets of the  
 95 form  $\mathcal{D} = \{\mathbf{x}^{(i)}, \boldsymbol{\theta}^{(i)}\}_{i=1}^n$  and use a deep generative model to learn a probabilistic mapping from data to parameters. Thus,  
 96 during the inference phase, the model generates samples  $\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}, \dots, \boldsymbol{\theta}^{(L)}$  from the observed data  $\mathbf{x}_{obs}$ , approximating the  
 97 target posterior  $p(\boldsymbol{\theta}|\mathbf{x} = \mathbf{x}_{obs})$ .

98 In the current work, we propose to implement and use a conditional invertible neural network (cINN) architecture. Previously,  
 99 INNs have been successfully employed to model data from astrophysics and medicine (23). We adapt the model to suit the task  
 100 of parameter estimation in the context of mathematical modeling (see Figure 1 for a full graphical illustration of the method)  
 101 and develop a reusable probabilistic architecture for full Bayesian likelihood-free inference on complex mathematical models.

**The Affine Coupling Block.** The basic building block of a cINN is the affine coupling block (ACB, see Figure 2a) (23, 24, 26).  
 Each cACB consists of four separate fully connected neural networks denoted as  $s_1(\cdot), s_2(\cdot), t_1(\cdot), t_2(\cdot)$ . An ACB is specifically  
 designed to be invertible, which means that in addition to a parametric mapping  $f_{\phi_{inv}} : \mathbb{R}^d \rightarrow \mathbb{R}^d$  it also learns the inverse  
 mapping  $f_{\phi_{inv}}^{-1} : \mathbb{R}^d \rightarrow \mathbb{R}^d$  "for free". Denoting the input vector of  $f_{\phi_{inv}}$  as  $\mathbf{u}$  and the output vector as  $\mathbf{v}$ , it follows that  
 $f(\mathbf{u}; \phi_{inv}) = \mathbf{v}$  and  $f^{-1}(\mathbf{v}; \phi_{inv}) = \mathbf{u}$ . Invertibility is achieved by splitting the input vector into two parts  $\mathbf{u} = (\mathbf{u}_1, \mathbf{u}_2)$  and  
 performing the following operations on the split input:

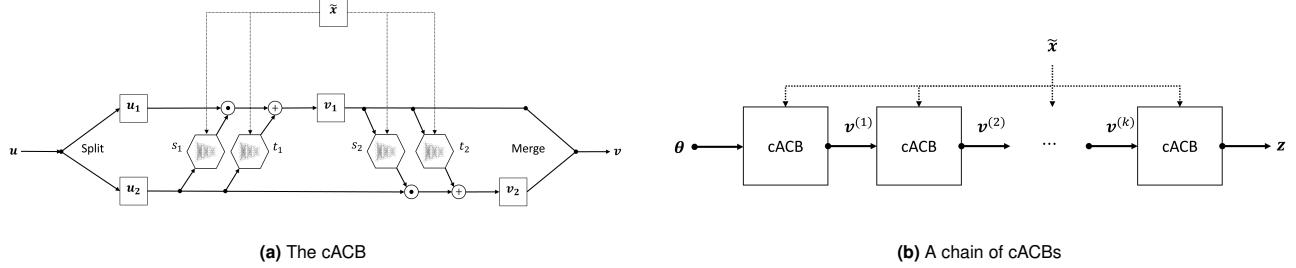
$$\mathbf{v}_1 = \mathbf{u}_1 \odot \exp(s_1(\mathbf{u}_2)) + t_1(\mathbf{u}_2) \quad [2]$$

$$\mathbf{v}_2 = \mathbf{u}_2 \odot \exp(s_2(\mathbf{v}_1)) + t_2(\mathbf{v}_1) \quad [3]$$

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:

$$\mathbf{u}_2 = (\mathbf{v}_2 - t_2(\mathbf{v}_1)) \odot \exp(-s_2(\mathbf{v}_1)) \quad [4]$$

$$\mathbf{u}_1 = (\mathbf{v}_1 - t_1(\mathbf{u}_2)) \odot \exp(-s_1(\mathbf{u}_2)) \quad [5]$$



**Fig. 2.** A diagram of the conditional version of the affine coupling block (cACB). **(a)** Each cACB consists of four internal networks performing the invertible operations described in the text; **(b)** In practice, we chain multiple cACBs to obtain higher representational capacity. Each cACB layer uses a fixed permutation to ensure that information about each parameter is encoded in each latent dimension of  $z$ .

An additional property of this design, which becomes relevant later for optimization, is that the operations of the ACB have tractable, and cheaply computable Jacobians (strictly upper or lower triangular matrices). Furthermore, the internal networks  $s_1(\cdot), s_2(\cdot), t_1(\cdot), t_2(\cdot)$  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 ACB. To ensure that the model is powerful enough to represent complicated distributions, we chain multiple ACBs, so that the output of each ACB becomes the input of the next (see Figure 2b). 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  $\phi_{inv}$ .

In our applications, the input to the first ACB is the parameter vector  $\theta$ , and the output of the final ACB, denoted hitherto as  $z$ , is encouraged to follow a  $d$ -dimensional spherical Gaussian via optimization (described in detail later), that is,  $p(z) = \mathcal{N}_d(z|\mathbf{0}, \mathbf{I})$ . Fixed permutation matrices are used before each ACB to ensure that each axis of the latent space encodes information from all components of  $\theta$ . In order to take into account the observed or simulated data  $x$ , each of the internal networks of each ACB is augmented to take  $x$  as an additional input -  $s_1(\cdot, x), s_2(\cdot, x), t_1(\cdot, x), t_2(\cdot, x)$  - so a complete pass through the entire invertible chain can be expressed as:

$$f(\theta; x, \phi_{inv}) = z \quad [6]$$

together with the inverse operation:

$$f^{-1}(z; x, \phi_{inv}) = \theta \quad [7]$$

This process can be interpreted as follows: the forward pass maps data-generating parameters to  $z$ -space using conditional information of  $x$ , while the inverse pass maps data points from  $z$ -space to the data-generating parameters of interest using the same conditional information provided by the data. In the next section, we describe the optimization procedure used to match the outputs of  $f^{-1}(z; x, \phi_{inv})$  to the posterior  $p(\theta|x)$ .

**Summary Network.** Since in practice the conditioning data set  $x$  can have variable number of input points (e.g., trial sizes, time points) and exhibit various redundancies, the cINN can profit from some form of dimensionality reduction applied to the data. Ideally, we want to avoid 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  $\tilde{x} = h(x; \phi_{sum})$  and learn the parameters of the summary network  $h$  jointly with those of the cINN chain via backpropagation. Thus, the current method remains completely end-to-end and is capable of generalizing to data sets of variable input size and structure.

**Learning the Posterior.** The cINN learns to approximate the posterior of model parameters by optimizing a maximum likelihood (ML) criterion. Broadly speaking, the goal of ML estimation is to find a set of parameters which maximize the probability of the data under a parametric model. In our case, we are interested in maximizing the expectation over all possible neural network parameters with respect to the parameters of the mathematical model:

$$\phi^* = \underset{\phi}{\operatorname{argmax}} \mathbb{E}_{\theta \sim p(\theta|x)} [p(\phi|\theta, x)] \quad [8]$$

Applying Bayes' rule to the posterior over all neural network parameters we obtain:

$$p(\phi|\theta, x) \propto p(\theta|x, \phi)p(\phi) \quad [9]$$

Note, that by maximizing Eq.9 we are maximizing the posterior over model parameters of interest  $p(\theta|x, \phi)$ . Thus, it remains to find a tractable expression for Eq.8 to be minimized by backpropagation given a finite number of simulated samples from the model. To this end, we recall that we can relate the pdf of  $\theta$  to that of  $z$  via the change of variable theorem:

$$p(\theta|x, \phi) = p(z) \left| \det \left( \frac{\partial z}{\partial \theta} \right) \right| \quad [10]$$

$$= p(f(\theta; x, \phi)) \left| \det \left( \frac{\partial f}{\partial \theta} \right) \right| \quad [11]$$

where  $\partial f / \partial \theta = \mathbf{J}_f$  is the Jacobian of the learned transformation  $f(\theta; \mathbf{x}, \phi)$  with respect to the input. Both terms in Eq. 11 are now tractable, since we have previously defined  $\mathbf{z}$  as following a spherical unit Gaussian, that is,  $p(\mathbf{z}) = (2\pi)^{-d/2} \exp(-\|\mathbf{z}\|_2^2)$  and the determinant of the Jacobian is easily computed as  $s_1(\mathbf{u}_2, \mathbf{x}) + s_2(\mathbf{v}_1, \mathbf{x})$  due to eqs. 2 and 3. We can now formulate the ML loss as the Monte-Carlo approximation of the negative logarithm of Eq. 8 for a batch of size  $m$ :

$$\mathcal{L}(\phi) = -\frac{1}{m} \sum_{i=1}^m \log(p(\phi | \theta^{(i)}, \mathbf{x}^{(i)})) \quad [12]$$

$$= -\frac{1}{m} \sum_{i=1}^m \log(p(\theta^{(i)} | \mathbf{x}^{(i)}, \phi)p(\phi)) \quad [13]$$

$$= -\frac{1}{m} \sum_{i=1}^m \left( \frac{\|f(\theta^{(i)}; \mathbf{x}^{(i)}, \phi)\|_2^2}{2} + \log \left| \det \left( \mathbf{J}_f^{(i)} \right) \right| \right) + \tau \|\phi\|_2^2 \quad [14]$$

where we place a Gaussian prior over the neural network parameters with  $\tau \equiv 1/\sigma^2$ , corresponding to a standard L2-regularization.

Minimizing Eq. 14 can be interpreted as searching for the optimal neural network parameters  $\phi^*$  which maximize the probability of model parameters  $\theta$  given data  $\mathbf{x}$ . This is exactly the probability we are concerned with in Bayesian inference. Note that our formulation maximizes the posterior of model parameters directly, in contrast to variational methods which optimize a lower bound on the posterior (22, 33). Once the backpropagation algorithm has settled to a local minimum of the ML loss, one can easily obtain samples from the approximate posterior  $p(\theta | \mathbf{x} = \mathbf{x}_{obs}, \phi = \phi^*)$ , based on an observed dataset  $\mathbf{x}_{obs}$ , by repeatedly sampling  $\mathbf{z}^{(l)} \sim \mathcal{N}_d(\mathbf{0}, \mathbf{I})$ , and then passing  $\mathbf{z}^{(l)}$  in reverse to the cINN in order to compute  $\theta^{(l)} = f^{-1}(\mathbf{z}^{(l)}; \mathbf{x}_{obs}, \phi^*)$  for  $l = 1, \dots, L$ . Figure 1b illustrates the inference phase of the method. It is worth noting that sampling a large number of parameter values from the approximate posterior takes a negligible amount of time, since it only requires a single pass through the cINN in reverse.

**Training the Networks.** We train all cINNs 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 set the weight regularization parameter  $\tau$  to a value of  $10^{-5}$ . With these settings, we perform 50 000 to 100 000 iterations (network updates) for each of the examples in this paper, and report the results on the trained network. Note, that we did not perform an extensive search for optimal values of the cINN hyperparameter, but use a cINN with 10 ACBs all examples in this paper (see SI for details of the cINN). All networks were implemented in Python using the *TensorFlow* library (32) and trained on a single-GPU machine equipped with NVIDIA® GTX1060 graphics card. Regarding the data generation step, we use two training approaches.

The first approach follows the classical approximate Bayesian computation approach to create a large *reference table* or grid of the form  $\mathbf{D} = \{\mathbf{x}^{(i)}, \theta^{(i)}\}_{i=1}^n$ . The reference table is then used as training data for the neural network and training continues for a pre-specified number of epochs through the entire reference table. A separate validation dataset is eventually used to assess the performance of the network. Training on large pre-simulated datasets separates the simulation from the training phase but can cause large memory overhead, as the reference table must be stored on disk and then loaded in chunks or in its entirety into memory.

The second approach follows a different strategy, which resembles ideas from *active learning* (34). Correspondingly, a dataset, or a batch of datasets, is created on the fly and then passed through the neural network. This training regime 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 switch to the first regime, since simulation and training in the active learning regime are tightly intertwined.

**Putting It All Together.** On an abstract level, our method requires three key ingredients: 1) a mathematical process model  $q(\theta)$  capable of simulating data  $\mathbf{x}$ ; 2) a prior distribution over the model parameters  $p(\theta)$  encoding our prior beliefs about plausible parameter values; 3) an invertible neural network  $f_\phi$  capable of approximating a large enough family of probability distributions (see Figure 2). In practice, a chain of up to 10 ACBs should suffice to learn most distributions encountered in the cognitive or life sciences, since they tend to be unimodal and relatively simple (in contrast to the distributions required to represent images or words in a spoken language). From these three ingredients, a universal and reusable sampler can be designed for likelihood-free Bayesian estimation of both tractable and intractable mathematical models. **Algorithm 1** describes the essential steps of the method using an arbitrary summary network and employing the active learning training regime.

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 local minimum. We experienced no instability or convergence issues during training with the ML loss. Note, that steps 12 – 15 of **Algorithm 1** can be executed in parallel with GPU support.

In what follows, we apply the method to a toy Bayesian regression example with conjugate priors, and then use it to estimate the parameters of challenging models from population dynamics, cognitive science, epidemiology, and genetics. Code for reproducing the results on all following examples is freely available at: <https://github.com/stefanradev93/cINN>.

---

**Algorithm 1** Bayesian likelihood-free inference with invertible neural networks

---

```

1: Training (via active learning):
2: repeat
3:   Sample a batch of  $\{\boldsymbol{\theta}^{(i)}\}_{i=1}^m$  from prior  $p(\boldsymbol{\theta})$ 
4:   Simulate a batch of datasets  $\{\mathbf{x}^{(i)}\}_{i=1}^m$  by running  $\mathbf{x}^{(i)} = q(\boldsymbol{\theta}^{(i)})$  for  $i = 1, \dots, m$ 
5:   Pass  $\{\mathbf{x}^{(i)}\}_{i=1}^m$  through summary network  $h(\mathbf{x}^{(i)}; \boldsymbol{\phi}_{sum})$  to obtain  $\{\tilde{\mathbf{x}}^{(i)}\}_{i=1}^m$ 
6:   Pass  $\{\boldsymbol{\theta}^{(i)}\}_{i=1}^m$  and  $\{\tilde{\mathbf{x}}^{(i)}\}_{i=1}^m$  through cINN  $f(\boldsymbol{\theta}^{(i)}; \mathbf{x}^{(i)}, \boldsymbol{\phi}_{inv})$  to obtain  $\{\mathbf{z}^{(i)}\}_{i=1}^m$ 
7:   Compute ML loss  $\mathcal{L}(\boldsymbol{\phi})$  according to Eq. 14
8:   Update neural network parameters  $\boldsymbol{\phi}$  via backpropagation
9: until convergence to  $\boldsymbol{\phi}^*$ 
10: Inference (given observed or test data  $\mathbf{x}_{obs}$ ):
11: Summarize the observed data by computing  $\tilde{\mathbf{x}}_{obs} = h(\mathbf{x}_{obs}, \boldsymbol{\phi}_{sum}^*)$ 
12: for  $l = 1, \dots, L$  do
13:   Sample  $\mathbf{z}^{(l)} \sim \mathcal{N}_d(\mathbf{0}, \mathbf{I})$ 
14:   Compute  $\boldsymbol{\theta}^{(l)} = f^{-1}(\mathbf{z}^{(l)}; \tilde{\mathbf{x}}_{obs}, \boldsymbol{\phi}_{inv})$ 
15: end for
16: Use  $\{\boldsymbol{\theta}^{(l)}\}_{l=1}^L$  to approximate the posterior  $p(\boldsymbol{\theta}|\mathbf{x}_{obs})$ 

```

---

166 **Results**

167 To assess the performance of our method on the following application examples, we consider a number of different metrics.  
168 To assess the precision of point estimates (posterior means), we compute the normalized root mean squared error (NRMSE)  
169 and the coefficient of determination ( $R^2$ ) between estimated and true parameter values. To assess the recovery of the full  
170 posterior, we compute the Kullback-Leibler (KL) divergence (35) between the true and the approximate distributions for the  
171 toy example, and use simulation-based calibration (SBC, (36)) for the other examples where the analytic posterior is not  
172 available in closed-form. Details for computing the NRMSE,  $R^2$ , and SCB can be found in the **SI**.

**Toy Example – Bayesian Regression.** As a proof-of-concept, we demonstrate the utility of our method in recovering the true analytic posteriors of the regression coefficients of a conjugate Bayesian regression model. To set the stage, assume we have observed a dataset  $\mathbf{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$  with  $\mathbf{x} \in \mathbb{R}^d$  and  $y \in \mathbb{R}$ . We stack all  $\mathbf{x}^{(i)}$  row-wise in a design matrix  $\mathbf{X}$  and model the vector of outcomes  $\mathbf{y}$  as being conditionally Gaussian given  $\mathbf{X}$ , that is,  $\mathbf{y} \sim \mathcal{N}_n(\mathbf{X}\boldsymbol{\theta}, a^{-1}\mathbf{I})$  where  $\boldsymbol{\theta} \in \mathbb{R}^d$  and  $a$  is the precision (inverse noise variance,  $a \equiv 1/\sigma_y^2$ ). We place a  $d$ -dimensional diagonal Gaussian prior on the regression coefficients centered at 0:  $\boldsymbol{\theta} \sim \mathcal{N}_d(\mathbf{0}, b^{-1}\mathbf{I})$  where  $b$  is the precision of the prior ( $b \equiv 1/\sigma_\theta^2$ ). Thus, the likelihood  $p(\mathbf{D}|\boldsymbol{\theta})$  admits the following proportionality:

$$p(\mathbf{D}|\boldsymbol{\theta}) \propto \exp\left(-\frac{a}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})\right) \quad [15]$$

Since the prior of  $\boldsymbol{\theta}$  is conjugate to the likelihood (both are Gaussian distributions), the posterior of  $\boldsymbol{\theta}$  is also Gaussian, given by:

$$p(\boldsymbol{\theta}|\mathbf{D}) \propto \exp\left(-\frac{a}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) - \frac{b}{2}\boldsymbol{\theta}^T\boldsymbol{\theta}\right) \quad [16]$$

Therefore, the posterior has the form  $p(\boldsymbol{\theta}|\mathbf{D}) = \mathcal{N}_d(\boldsymbol{\theta}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$  where  $\boldsymbol{\Lambda}$  denotes the posterior precision matrix (inverse covariance matrix),  $\boldsymbol{\mu}$  the posterior mean vector, which are computed as follows:

$$\boldsymbol{\Lambda} = a\mathbf{X}^T\mathbf{X} + b\mathbf{I} \quad [17]$$

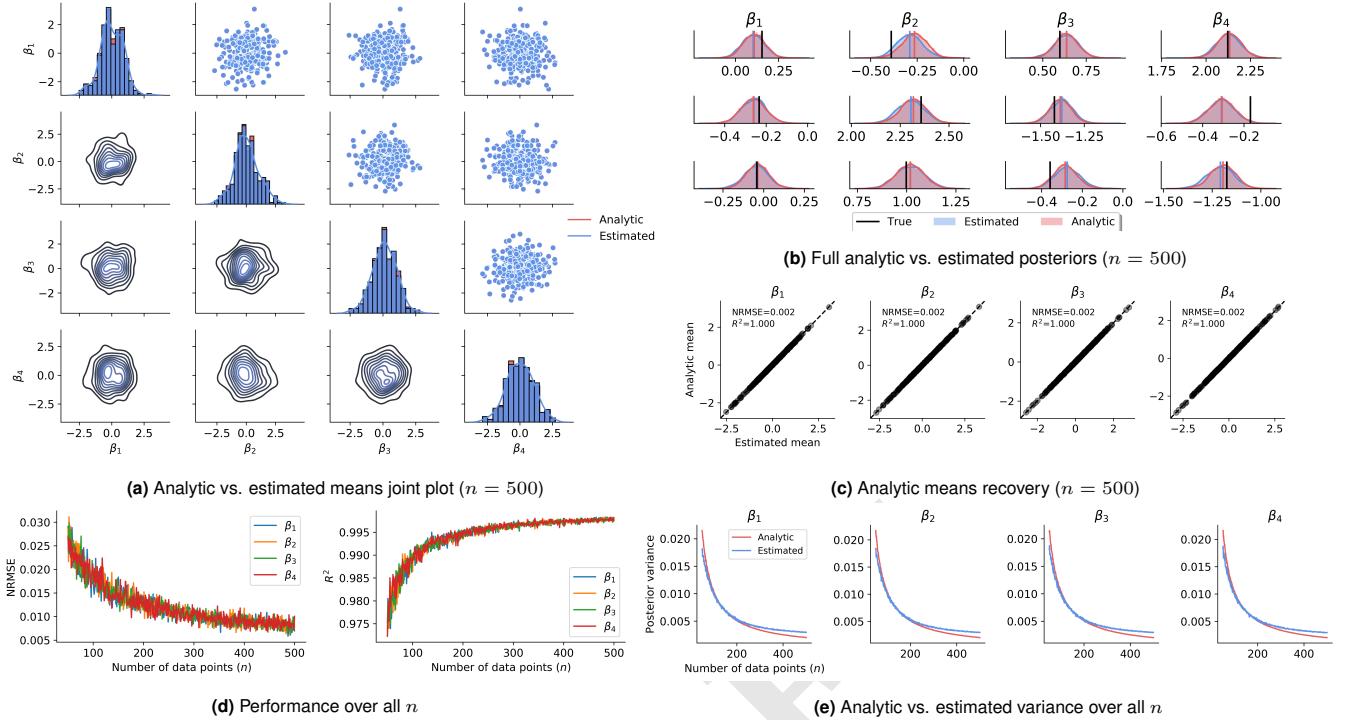
$$\boldsymbol{\mu} = a\boldsymbol{\Lambda}^{-1}\mathbf{X}^T\mathbf{y} \quad [18]$$

173 Thus, for known  $a$  and  $b$ , the posterior can be easily computed. Even though in real-world applications  $a$  is usually not known  
174 and a hierarchical model is used instead, the current example is good for testing the utility of our method.

175 For the following application, we set  $d = 4$ , and  $a = b = 1$ . The design matrices for each iteration contain a variable number  
176  $n$  of *i.i.d.* data points drawn from a unit Gaussian  $\mathbf{x}^{(i)} \sim \mathcal{N}_4(\mathbf{0}, \mathbf{I})$  for  $i = 1, \dots, n$ . The number of trials is drawn from a  
177 uniform distribution  $n \sim \mathcal{U}(50, 500)$  at each training iteration (Lines 2-9 of **Algorithm 1**). Training the networks took approx.  
178 half a day with the active learning approach. Inference on 1000 datasets with 2000 posterior samples per parameter took  
179 approx. 5.5 seconds.

180 The results on the toy Bayesian regression are depicted in **Figure 3**. The approximate posterior means show negligible  
181 deviations from the analytic posterior means as quantified by very small NRMSEs (as small as 0.002) and very high  $R^2$  (as  
182 high as 1.0) over all  $n$  sampled during training. This suggests near-perfect estimation of the true posterior means. Further, the  
183 estimates become increasingly more accurate, as the number of data points  $n$  increases. An inspection of the posterior variances  
184 over the different  $n$  reveals that the estimated variance follows closely the decrease in analytic variance with increasing  $n$ .  
185 However, the analytic variance is slightly underestimated at smaller  $n$  and slightly overestimated at larger  $n$ . This pattern is

also revealed by the KL divergence plot (see **SI**). This result might be attributable to an underexpressive summary network; another possibility is that the networks need to be trained longer with smaller learning rate decay.



**Fig. 3.** Results on the Bayesian toy regression example. **(a)** Pair-plots of the analytic vs. estimated posterior means on the whole validation sample. We observe a near-perfect overlap between analytic and estimated posterior means and no spurious covariances; **(b)** Some example draws from the estimated posteriors. Visual inspections reveals a close match between analytic and estimated posteriors; **(c)** Posterior mean recovery is almost perfect at  $n = 500$ ; **(d)** NRMSE and  $R^2$  over all  $n$ . Performance improves with increasing number of data points; **(e)** Analytic vs. estimated variance over all  $n$ . The analytic variance is slightly underestimated at lower  $n$  and slightly overestimated at higher  $n$ .

**Example 1 - The Ricker Model.** Discrete population dynamics models describe how the number of individuals in a population changes over discrete units of time (7). In particular, the Ricker model describes the number of individuals  $x_{t+1}$  in generation  $t+1$  as a function of the number of individuals in the previous generation  $t$  by the following non-linear equation:

$$x_t \sim Pois(\rho N_t) \quad [19]$$

$$N_{t+1} = r N_t e^{-N_t + \epsilon_t} \quad [20]$$

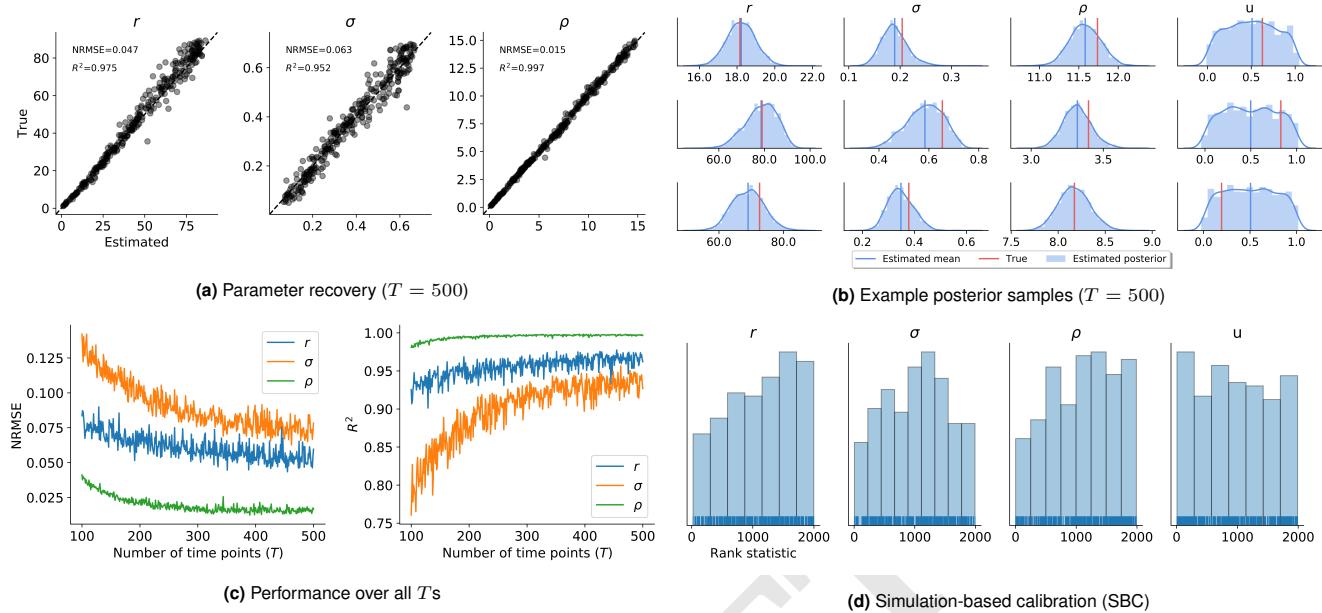
for  $t = 1, \dots, T$  where  $N_t$  is the expected number of individuals at time  $t$ ,  $r$  is the growth rate,  $\rho$  is a scaling parameter and  $\epsilon_t \sim \mathcal{N}(0, \sigma^2)$  is random Gaussian noise (19). The likelihood function for the Ricker model is not available in closed-form, and the model is known to exhibit chaotic behavior. Thus, it is a suitable candidate for likelihood-free inference. The parameter estimation task is thus to recover  $\theta = (\rho, r, \sigma)$  from the observed one-dimensional time-series data  $\mathbf{x} = (x_1, x_2, \dots, x_T)$  where each  $x_t \in \mathbb{N}$ .

During training of the networks, 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 (see the **SI** for more details about the simulation). Training the networks took approx. half a day with the active learning approach. Inference on 1000 datasets with 2000 posterior samples per parameter took approx. 6 seconds.

What if the data does not contain 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  $\theta$  and train the model with this additional dummy parameter. We expect that the networks learn to ignore this dummy parameter, that is, the estimated posterior of  $u$  is (approximately) equal to the uniform prior.

The results on the Ricker model are depicted in Figure 4. As evident from the graphs, parameter recovery becomes better when more time points with data are available (Figure 4c). At  $T = 500$ , the NRMSEs range between 0.015 and 0.063, and the  $R^2$  metrics between 0.997 and 0.952, indicating very good recovery of the posterior means (Figure 4a). The parameter  $\sigma$  seems to be the hardest to recover. Inspecting the full posteriors, we further see that the posterior distribution of the dummy noise variable  $u$  closely resembles the prior, as expected due to the complete lack of mutual information between the data  $x$  and  $u$  (Figure 4b). Finally, the plots of the rank statistic computed for SCB suggest no systematic distortions of the posterior across all parameters (Figure 4d). Interpreting deviations from uniformity according to (36), the approximate posteriors of  $r$  and  $\rho$  slightly underestimate the true posterior means, whereas the approximate posterior of  $\sigma$  tends to overestimate the true

209 posterior variance. These deviations appear to be due to the fact that recovery worsens at extreme values of the parameters.  
 210 This is unsurprising, as the data generated with these parameters is highly implausible, which in some cases might even render  
 211 a model unidentifiable.



**Fig. 4.** Results on the Ricker model. **(a)** Parameter recovery for the maximum number of generations used during training ( $T = 500$ ); **(b)** Example posteriors for three test datasets. We observe that the posterior of the uniform noise variable  $u$  is equal to the prior, i.e., the method detects that no information is present in data for this variable; **(c)** NRMSE and  $R^2$  performance metrics over all  $T$ 's used in training. We observe that the recovery remains good over all  $T$ 's, and becomes progressively better as more data is available; **(d)** Plots of the rank statistics indicative of the accuracy of the full posterior. Accordingly, the approximate posteriors of  $r$  and  $\rho$  slightly underestimate the true posterior means, whereas the approximate posterior of  $\sigma$  tends to overestimate the true posterior variance.

**Example 2 - The Lévy-Flight Model.** Evidence accumulator models (EAMs) describe human decision making by a set of neurocognitively motivated parameters (37). EAMs are most often applied to choice reaction times (RT) data to obtain an estimate of the neurocognitive processes governing observed RT distributions in human (or animal) participants. Most EAM variants share four underlying assumptions: *i*) information about a stimulus (response option) is accumulated continuously through time; *ii*) stochasticity in the form of noisy accumulation ensures variability; *iii*) empirical response times can be decomposed into a decision time component and a non-decision time component accounting for pre-decisional perceptual (encoding time) and post-decisional motor processes (response execution); and *iv*) a decision is met when the activation of an accumulator reaches a threshold. In its most general formulation, the forward model of EAMs takes the form of a stochastic differential equation given by (4):

$$dx = vdt + cd\xi \quad [21]$$

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

EAMs are particularly amenable for likelihood-free inference, since the likelihood of some members of this model family turn out to be intractable (38). This intractability has precluded many interesting applications and empirically driven model refinements. Here, we apply our method to estimate the parameters of the recently proposed Lévy-Flight Model (LFM, (39)). The LFM assumes an *alpha-stable* noise distribution of the evidence accumulation process in order to model "jumps" in the decision process. However, the inclusion of *alpha-stable* noise leads to a model with intractable likelihood; further, to our knowledge, a fully Bayesian treatment of the model is still missing in the literature. The forward equation of the LFM is given by:

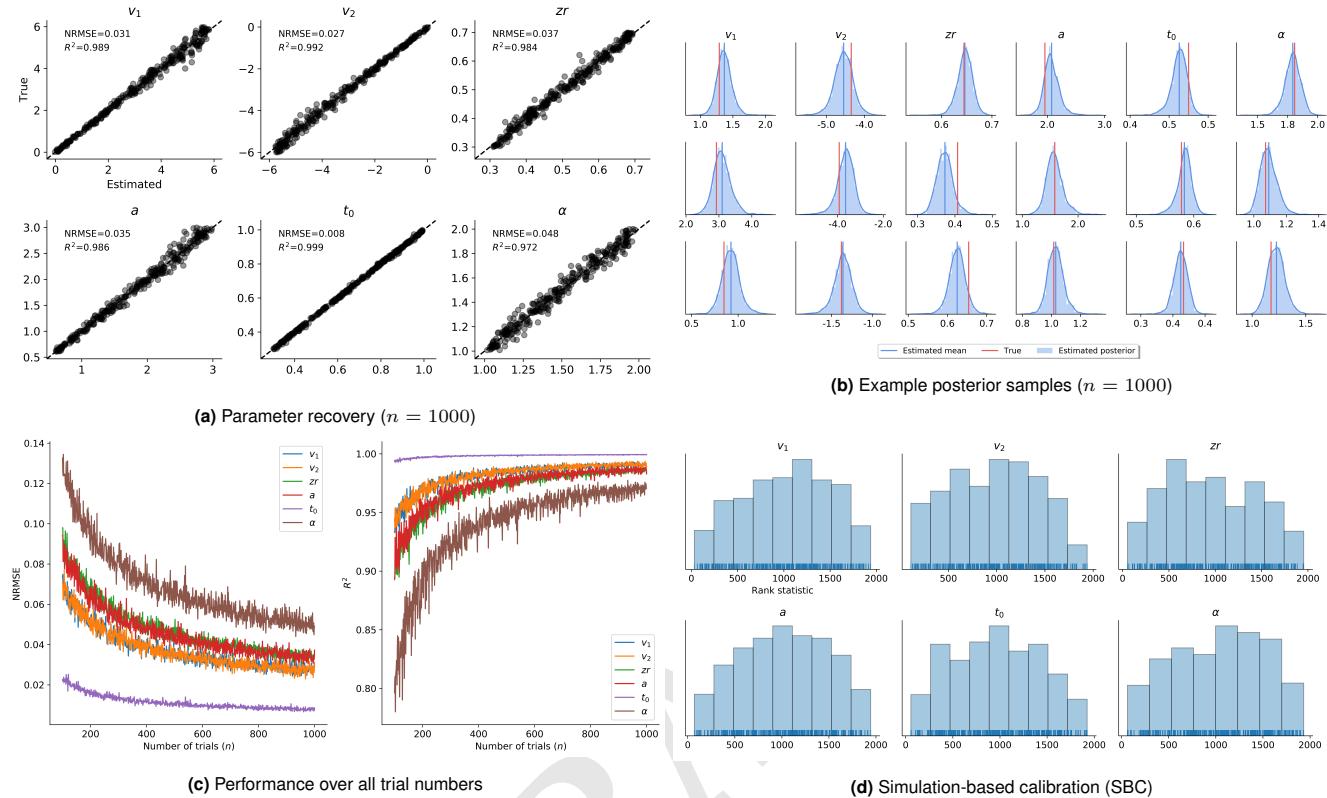
$$dx = vdt + \xi dt^{1/\alpha} \quad [22]$$

$$\xi \sim \text{AlphaStable}(\alpha, 0, 1, 0) \quad [23]$$

215 The LFM has three additional parameters: the threshold  $a$  determining the amount of evidence needed for the termination of a  
 216 decision process; a relative starting point,  $zr$ , determining the amount of starting evidence available to the accumulator before  
 217 the actual decision alternatives are presented; and an additive non-decision time  $t_0$ .

218 During training of the networks, we simulate response times data from two experimental conditions with two different  
 219 drift rates (see SI for details of the simulation), since this 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  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$  containing variable number or RT trials. The number of trials is drawn from a uniform distribution  $n \sim \mathcal{U}(100, 1000)$  at each training iteration. Training the networks took a little less than a day with the active learning approach. Inference on 1000 datasets with 2000 posterior samples per parameter took approx. 7.39 seconds.



**Fig. 5.** Results on the LFM model. **(a)** Parameter recovery for the maximum number of trials used during training ( $n = 1000$ ); **(b)** Example posteriors for three test datasets; **(c)** NRMSE and  $R^2$  performance metrics over all  $n$  trials used during training. Again, we observe that the recovery remains overall very good, and becomes progressively better as more data is available; **(d)** Plots of the rank statistics indicative of the accuracy of the full posterior. Accordingly, the approximate posteriors tend to overestimate the true posterior variance.

The results on the LFM model are depicted in [Figure 5](#). To our knowledge, this is the first Bayesian treatment of the LFM, as its intractability makes traditional methods prohibitively slow. We observe excellent recovery of all LFM parameters with NRMSEs ranging between 0.008 and 0.048 and  $R^2$  between 0.972 and 0.999. Further, estimation remains very good across all trial sizes, with goodness increasing as more trials become available. The parameter  $\alpha$  appears to be most challenging to estimate, requiring more data for good estimation, whereas the non-decision time parameter  $t_0$  is almost perfectly reconstructed for all trial sizes. Last, inspecting the SCB plots, we notice that the histograms tend to be slightly peaked, indicating a slight overestimation of the posterior variance ([36](#)).

**Example 3 – The Stochastic SIR Model.** Compartmental models in epidemiology are used to describe the stochastic dynamics of infectious diseases as they spread over a population of individuals ([9](#), [10](#), [40](#)). The parameters of compartmental models encode important characteristics of diseases, such as the rates of infection or recovery from the disease. The stochastic SIR model describes the transition dynamics of  $N$  individuals between three discrete states: susceptible ( $S$ ), infected ( $I$ ), and recovered ( $R$ ). The transition dynamics are given by the following equations:

$$\Delta S = -\Delta N_{SI} \quad [24]$$

$$\Delta I = \Delta N_{SI} - \Delta N_{IR} \quad [25]$$

$$\Delta R = \Delta N_{IR} \quad [26]$$

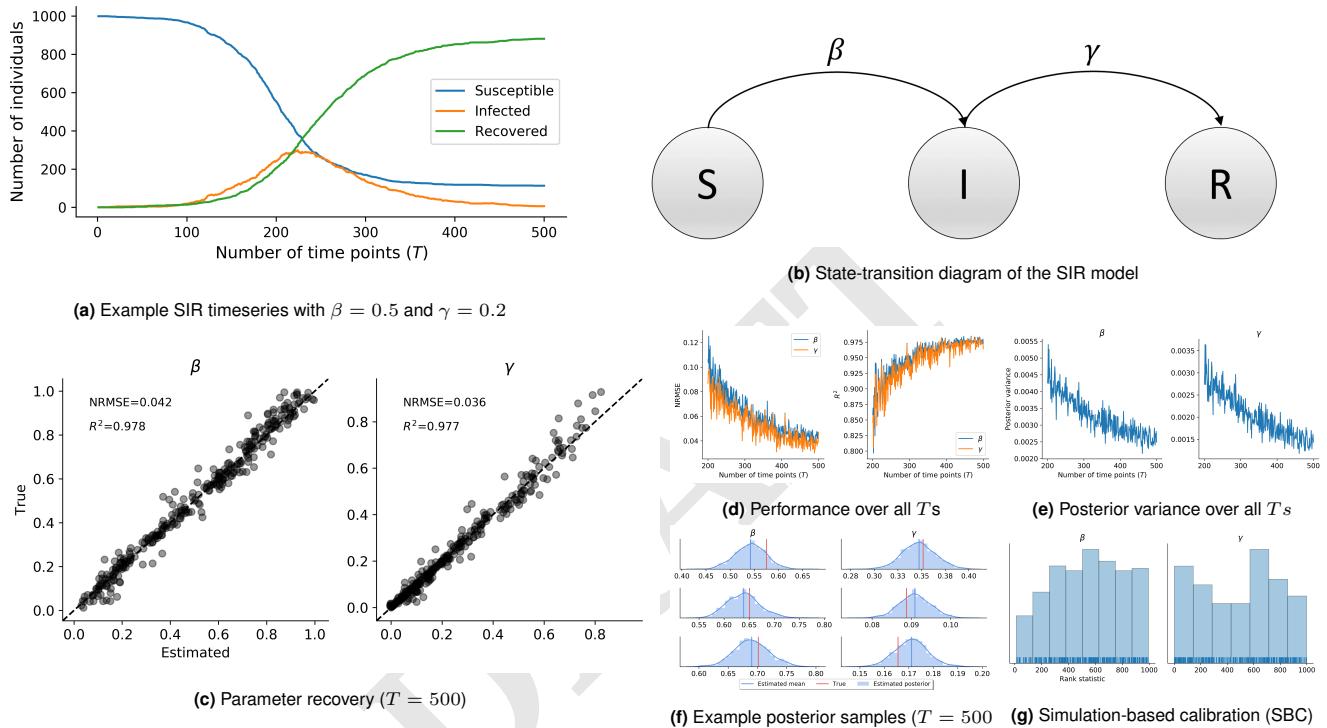
$$\Delta N_{SI} \sim \text{Binomial}(S, 1 - \exp(-\beta \frac{I}{N} \Delta t)) \quad [27]$$

$$\Delta N_{IR} \sim \text{Binomial}(I, 1 - \exp(-\gamma \Delta t)) \quad [28]$$

where  $S + I + R = N$  give the number of susceptible, infected, and recovered individuals, respectively. The parameter  $\beta$  controls the transition rate from being susceptible to infected, and  $\gamma$  controls the transition rate from being infected to recovered (see

233 **Figure 6**). The number of individuals moving from  $S$  to  $I$ , given by  $\Delta N_{SI}$ , and the number of people moving from  $I$  to  $R$ ,  
 234 given by  $\Delta N_{IR}$ , over a time interval  $\Delta t$  are modeled as binomial random variables. The above listed stochastic system has no  
 235 known analytic solution and thus requires numerical simulation methods for recovering parameter values from data. Cast as a  
 236 parameter estimation task, the challenge is to recover  $\theta = \{\beta, \gamma\}$  from three dimensional time-series data  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T)$   
 237 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  
 238 at time  $t$ .

239 During training of the networks, we simulate time-series from the stochastic SIR model with varying lengths. The number  
 240 of time points  $T$  is drawn from a uniform distribution  $T \sim \mathcal{U}(200, 500)$  at each training iteration (see the **SI** for more details  
 241 about the simulation). Usually, at lower  $T$ s, the system has not converged to an equilibrium (i.e., not all individuals have  
 242 transitioned from being  $I$  to  $R$ ). Thus, it is especially interesting to see if our method can recover the parameters, even if  
 243 the process dynamics are still unfolding over time. Training the networks took approx. half a day with the active learning  
 244 approach. Inference on 1000 datasets with 2000 posterior samples per parameter took approx. 1.1 seconds.



**Fig. 6.** Results on the stochastic SIR model. **(a)** Example SIR timeseries generated with  $\beta = 0.5$  and  $\gamma = 0.2$ ; **(b)** State transition diagram of the SIR model; **(c)** Parameter recovery depicting also NRMSE and  $R^2$  metrics; **(d)** NRMSE and  $R^2$  performance over all  $T$ s seen by the networks during training. We observe similar patterns as in the previous examples; **(e)** Posterior variances of the SIR parameters over all  $T$ s; **(f)** Example full posteriors of the parameters recovered from three test datasets; **(g)** Plots of the rank statistic, according to which the posterior of  $\beta$  slightly overestimates the true mean.

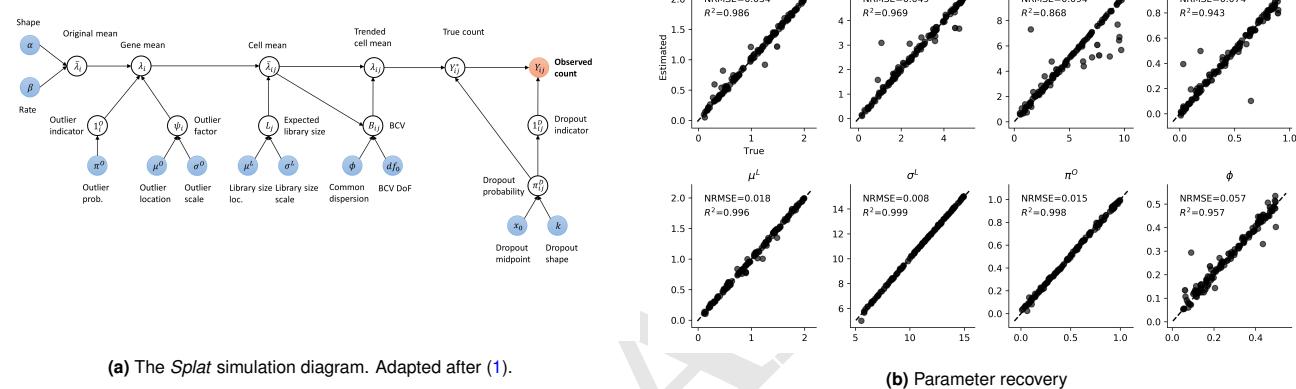
245 The results on the SIR model are depicted in **Figure 6**. Similar to the previous examples, we observe very good recovery of  
 246 the true parameters, with NRMSE at  $T = 500$  around 0.04, and  $R^2$ 's around 0.98. Further, we observe decent performance  
 247 even at smaller  $T$ s, with performance increasing as  $T$  increases. Conversely, the estimated posterior variance drops as  $T$   
 248 increases, capturing the fact that more information becomes available to the networks. Finally, the SCB plots indicate a slight  
 249 underestimation of the true posterior means by the estimated posterior means. This trend is visible when  $\beta$  and  $\gamma$  both assume  
 250 very similar values, and the generated timeseries are far from equilibrium at the final time point  $T$ . Thus, such datasets contain  
 251 very little information about the parameters given the considered time range of  $T_{max} = 500$ .

252 **Example 4 - Single-Cell RNA Sequencing.** Single-cell RNA sequencing (scRNA-seq) is a method to uncover the dynamics of  
 253 gene expression within single cells (1, 41). Simulation models for scRNA-seq attempt to mimic the process of sequencing  
 254 real data samples by combining statistical and algorithmic procedures. The recently developed *Splat* simulation model (1)  
 255 implements a hierarchical model where the mean expression levels for multiple genes are samples from a Gamma distribution,  
 256 and the number of times a gene has been sequenced in a cell is then sampled from a Poisson distribution. The output of the  
 257 *Splat* simulation is thus a matrix  $\mathbf{X}$  of *Gene x Cell* counts. **Figure 7a** depicts the *Splat* algorithm for simulating scRNA-seq data  
 258 (1). As a last example for our method, we use simulations from the *Splat* model and attempt to recover the data-generating  
 259 parameters. We simulate 500 000 *Gene x Cell* count matrices with different parameter settings (see **SI** for details of the  
 260 simulation) and discard implausible simulations (e.g., matrices where more than 50% of the counts are 0). The parameters  
 261 considered for estimation, along with a short description, are listed in **Table 1**. The parameter estimation task is thus to recover

**Table 1. Splat target parameters estimated by our method**

Parameter	Name	Description
$\alpha$	Mean shape	Shape of the mean gene expression gamma distribution
$\beta$	Mean rate	Rate of the mean gene expression gamma distribution
$\mu^L$	Library size location	Location of the library size log-normal distribution
$\sigma^L$	Library size scale	Scale of the library size log-normal distribution
$\pi^O$	Outlier probability	Probability that a gene is an expression outlier
$\mu^O$	Outlier location	Location of the expression outlier factor log-normal distribution
$\sigma^O$	Outlier scale	Scale of the expression outlier factor log-normal distribution
$\phi$	Common BCV	Common BCV dispersion across all genes

262  $\theta = (\alpha, \beta, \mu^L, \sigma^L, \pi^O, \mu^O, \sigma^O, \phi)$  from an observed count matrix  $\mathbf{X}$ . We train the networks for 50 epochs through the entire  
263 dataset and evaluate the performance on a separate validation set of 500 count matrices.



**Fig. 7.** Results on the Splat scRNA simulation model. **(a)** The Splat simulation diagram; **(b)** Parameter recovery.

264 The results on the scRNA example are depicted in Figure 7. We observe excellent parameter recovery on the validation set.  
265 **TODO: Train the full method. These results are preliminary.**

## Discussion

266 In the current work, we proposed and explored a novel end-to-end likelihood-free method which uses invertible neural networks  
267 to perform approximate Bayesian inference on any mathematical process model. We demonstrated the utility of the method  
268 by applying it to models from various scientific domains dealing with several different data formats and data-generating  
269 mechanisms. Further, we explored two possible training approaches suitable for different simulation scenarios, namely an active  
270 learning approach, and a reference-table approach. Both training approaches lead to excellent recovery of the true parameters  
271 throughout the examples considered in the current work.

272 Our method combines the universal approximation power of deep learning methods with the crucial uncertainty quantification  
273 assets of Bayesian inference (30, 31). Besides being capable of performing full Bayesian inference on intractable mathematical  
274 models, our method provides a general framework for designing reusable *parameter estimation machines* for various research  
275 domains. Moreover, it can also prove as a viable alternative in modeling contexts where standard inference methods are  
276 available, but inference is nevertheless computationally inconvenient.

277 Inspired by previous machine learning approaches to likelihood-free inference (5, 6, 18–22), our method shares many of the  
278 advantages of these methods and further overcomes some of their important limitations.

279 First, the introduction of separate summary and inference neural network modules renders the invertible inference module  
280 independent of the shape or the size of the observed data. This is achieved by learning a fixed-size vector representation of the  
281 data through the summary module in an automatic, data-driven manner. This is particularly useful in settings where the  
282 most informative summary statistics are not known and thus potential information is lost through the choice of suboptimal  
283 summary function. However, if informative or even *sufficient* statistics are available in a given domain, one might dispose with  
284 the summary module altogether and feed the summary statistics directly to the CINN.

285 Second, we showed that the ML loss directly maximizes the posterior over model parameters of interest, which is in contrast  
286 to ELBO-based methods which optimize a lower-bound on the posterior (22, 33). Thus, inference is exact when the ML is  
287 minimized (24, 26). This *optimal performance* claim is confirmed by our toy Bayesian regression example, in which we observe  
288 negligibly small deviations from the form of the true posterior. Further, researchers are often interested in some summary

of the posterior, for instance, the posterior mean or the posterior variance (18, 20). We demonstrated that our method exhibits excellent recovery of the posterior means throughout the examples. We also showed that parameter recovery becomes better with increasing number of observations, whereas the posterior variance of the parameters becomes smaller. These are important and highly desirable properties of any parameter estimation method, as it mirrors the increase in information following increasing number of data points.

Third, the largest computational cost of our method is paid during training. Once trained, the networks can be used and reused to perform inference on large numbers of datasets within seconds and across a given research domain. Indeed, there are many instances of research domains where a single model is extensively explored and independently fitted by multiple researches to test scientific hypotheses (1, 37, 39, 42). These research domains are expected to benefit the most from learning the “model universe” once and then inverting the model multiple times for rapid inference on multiple datasets. In this regard, our method is similar to the recently introduced prepaid method (19) which uses a database of pre-computed summary statistics and nearest-neighbors for inference. Note, however, that our method does not need to store training data, since the *knowledge* about the relationship between data and parameters is compressed into the networks’ weights. This not only makes the global sharing of pre-trained parameter estimation networks across researchers extremely easy, but also makes their local storage very efficient. Moreover, all computations involved in our method benefit from a high degree of parallelism and can thus utilize the advantages of modern GPUs.

These advantages notwithstanding, some limitations of the method deserve mention. Even though high-level deep learning libraries, such as *TensorFlow* or *Torch*, allow for rapid and relatively straightforward development of various neural network architectures, the implementational burden associated with the current method is still reasonably high. In order to ease the understanding and independent application of the method, we provide fully functioning code to reproduce and study all of the examples tackled in this paper (<https://github.com/stefanradev93/cINN>). In addition, we also provide implementation of all tools for performance validation used throughout the paper. Moreover, we are currently developing a general user-friendly software, which should abstract away most of the methodological complexities from the user. Another potential shortcoming of the method is the seemingly overwhelming number of hyperparameters that might require fine-tuning by the user for optimal performance on a given task. However, we observe that many of the default hyperparameter values are sufficient to achieve excellent performance, and starting with a relatively large default network of 10 ACBs does not appear to hurt performance or destabilize training, even if the mathematical model to be learned is relatively simple. We expect that a single architecture should be able to perform well on almost all models from a given domain (i.e., a single architecture for decision-making models). Future research should investigate the question of generality by applying the method to challenging parameter estimation tasks across different research domains.

**ACKNOWLEDGMENTS.** Please include your acknowledgments here, set in a single paragraph. Please do not include any acknowledgments in the Supporting Information, or anywhere else in the manuscript.

1. Zappia L, Phipson B, Oshlack A (2017) Splatter: simulation of single-cell rna sequencing data. *Genome biology* 18(1):174.
2. Beaumont MA, Zhang W, Balding DJ (2002) Approximate bayesian computation in population genetics. *Genetics* 162(4):2025–2035.
3. Palestro JJ, Sederberg PB, Osth AF, Van Zandt T, Turner BM (2018) *Likelihood-free methods for cognitive science*. (Springer).
4. Usher M, McClelland JL (2001) The time course of perceptual choice: the leaky, competing accumulator model. *Psychological review* 108(3):550.
5. Hwang SJ, Tao Z, Kim WH, Singh V (2018) Conditional recurrent flow: Conditional generation of longitudinal samples with applications to neuroimaging. *arXiv preprint arXiv:1811.09897*.
6. Lueckmann JM, et al. (2017) Flexible statistical inference for mechanistic models of neural dynamics in *Advances in Neural Information Processing Systems*. pp. 1289–1299.
7. Wood SN (2010) Statistical inference for noisy nonlinear ecological dynamic systems. *Nature* 466(7310):1102.
8. Geritz SA, Kisdi É (2004) On the mechanistic underpinning of discrete-time population models with complex dynamics. *Journal of Theoretical Biology* 228(2):261–269.
9. Keeling MJ, Rohani P (2011) *Modeling infectious diseases in humans and animals*. (Princeton University Press).
10. Hethcote HW (2000) The mathematics of infectious diseases. *SIAM review* 42(4):599–653.
11. Csillér K, Blum MG, Gaggiotti OE, François O (2010) Approximate bayesian computation (abc) in practice. *Trends in ecology & evolution* 25(7):410–418.
12. Toni T, Stumpf MP (2009) Simulation-based model selection for dynamical systems in systems and population biology. *Bioinformatics* 26(1):104–110.
13. Turner BM, Sederberg PB (2014) A generalized, likelihood-free method for posterior estimation. *Psychonomic bulletin & review* 21(2):227–250.
14. Sunnåker M, et al. (2013) Approximate bayesian computation. *PLoS computational biology* 9(1):e1002803.
15. Mertens UK, Voss A, Radev S (2018) Abrox—a user-friendly python module for approximate bayesian computation with a focus on model comparison. *PLoS one* 13(3):e0193981.
16. Frazier DT, Martin GM, Robert CP, Rousseau J (2018) Asymptotic properties of approximate bayesian computation. *Biometrika* 105(3):593–607.
17. Mertens UK (2019) Ph.D. thesis.
18. Radev ST, Mertens UK, Voss A, Köthe U (2019) Towards end-to-end likelihood-free inference with convolutional neural networks. *British Journal of Mathematical and Statistical Psychology*.
19. Mestdagh M, Verdonck S, Meers K, Loossens T, Tuerlinckx F (2018) Prepaid parameter estimation without likelihoods. *arXiv preprint arXiv:1812.09799*.
20. Raynal L, et al. (2018) Abc random forests for bayesian parameter inference. *Bioinformatics* 35(10):1720–1728.
21. Jiang B, Wu Ty, Zheng C, Wong WH (2017) Learning summary statistic for approximate bayesian computation via deep neural network. *Statistica Sinica* pp. 1595–1618.
22. Papamakarios G, Murray I (2016) Fast  $\epsilon$ -free inference of simulation models with bayesian conditional density estimation in *Advances in Neural Information Processing Systems*. pp. 1028–1036.
23. Ardizzone L, et al. (2018) Analyzing inverse problems with invertible neural networks. *arXiv preprint arXiv:1808.04730*.
24. Kingma DP, Dhariwal P (2018) Glow: Generative flow with invertible 1x1 convolutions in *Advances in Neural Information Processing Systems*. pp. 10215–10224.
25. Grover A, Dhar M, Ermon S (2018) Flow-gan: Combining maximum likelihood and adversarial learning in generative models in *Thirty-Second AAAI Conference on Artificial Intelligence*.
26. Dinh L, Sohl-Dickstein J, Bengio S (2016) Density estimation using real nvp. *arXiv preprint arXiv:1605.08803*.
27. Bloem-Reddy B, Teh YW (2019) Probabilistic symmetry and invariant neural networks. *arXiv preprint arXiv:1901.06082*.
28. Gers FA, Schmidhuber J, Cummins F (1999) Learning to forget: Continual prediction with lstm.
29. Goodfellow I, Bengio Y, Courville A (2016) *Deep learning*. (MIT press).
30. Kendall A, Gal Y (2017) What uncertainties do we need in bayesian deep learning for computer vision? in *Advances in neural information processing systems*. pp. 5574–5584.
31. Gelman A, et al. (2013) *Bayesian data analysis*. (Chapman and Hall/CRC).
32. Abadi M, et al. (2016) Tensorflow: A system for large-scale machine learning in 12th {USENIX} Symposium on Operating Systems Design and Implementation ({ OSDI } 16). pp. 265–283.
33. Kingma DP, Welling M (2014) Auto-encoding variational bayes. *stat* 1050:1.
34. Mnih V, et al. (2015) Human-level control through deep reinforcement learning. *Nature* 518(7540):529.
35. Hershey JR, Olsen PA (2007) Approximating the kullback leibler divergence between gaussian mixture models in 2007 IEEE International Conference on Acoustics, Speech and Signal Processing- ICASSP'07. (IEEE), Vol. 4, pp. IV–317.
36. Talts S, Betancourt M, Simpson D, Vehtari A, Gelman A (2018) Validating bayesian inference algorithms with simulation-based calibration. *arXiv preprint arXiv:1804.06788*.
37. Ratcliff R, McKoon G (2008) The diffusion decision model: theory and data for two-choice decision tasks. *Neural computation* 20(4):873–922.
38. Miletic S, Turner BM, Forstmann BU, van Maanen L (2017) Parameter recovery for the leaky competing accumulator model. *Journal of Mathematical Psychology* 76:25–50.

- 361 39. Voss A, Lerche V, Mertens U, Voss J (2019) Sequential sampling models with variable boundaries and non-normal noise: A comparison of six models. *Psychonomic bulletin & review* pp. 1–20.
- 362 40. Sahneh FD, Vajdi A, Shakeri H, Fan F, Scoglio C (2017) Gemfsim: a stochastic simulator for the generalized epidemic modeling framework. *Journal of computational science* 22:36–44.
- 363 41. Ozsolak F, Milos PM (2011) Rna sequencing: advances, challenges and opportunities. *Nature reviews genetics* 12(2):87.
- 364 42. De Valpine P, Hastings A (2002) Fitting population models incorporating process noise and observation error. *Ecological Monographs* 72(1):57–76.

DRAFT