

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

Stefan T. Radev¹, Ulf K. Mertens¹, Andreas Voss¹, Lynton Ardizzone², and Ullrich Köthe²

¹Institute of Psychology, Heidelberg University, Hauptstr. 47-51, 69117 Heidelberg, Germany; ²Heidelberg 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 June 29, 2019

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

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

Mathematical models are formal descriptions of scientific theories allowing a clear and unambiguous way to formulate and test scientific hypotheses about probabilistic phenomena in a probabilistic world. In its most abstract form, a mathematical model is specified by a set of parameters θ and a forward model q mimicking the process by which manifest data x arise from latent parameters:

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

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

7 Once a mathematical model has been formulated, the next step consists of fitting the model to experimental or observational
8 data and recovering the parameters of interest. However, estimating the parameters of a mathematical model from data
9 can quickly become one of the most tenacious challenges in applications to real-world problems. It is also one of the most
10 important ones to be tackled, since without reliable parameter estimation methods, it is impossible to test the utility of a
11 model, regardless of its sophistication or theoretical appeal. Idealized parameter estimation involves computing the inverse
12 (backward) model $\theta = q^{-1}(x)$ exactly. However, due to noise, inherent stochasticity or loss of information, the inverse usually
13 does not exist, so researchers need to resort to the sophisticated frameworks of Bayesian or frequentist inference.

14 However, as mathematical models and processes under description become increasingly complex, parameter estimation and
15 model selection can quickly become intractable with standard Bayesian and frequentist method. Complex models specified by
16 a generative stochastic mechanism do not always provide a closed-form solution for the *likelihood function* (3, 10, 11). This
17 poses great difficulties for Bayesian and frequentist methods alike, since both depend explicitly on the numerical evaluation
18 of a likelihood function as a proxy for assessing model fit to data. Even if a likelihood function is available in closed-form,

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")

²To whom correspondence should be addressed. E-mail: stefan.ralev@psychologie.uni-heidelberg.de

19 inference may be prohibitively slow for real-world applications. In this case, enforcing simplifying distributional assumptions
20 (i.e., independence or Gaussian assumptions) on the likelihood can increase the speed of inference, but can also lead to model
21 misspecifications and dramatically incorrect estimates. Therefore, the need for powerful and reliable likelihood-free estimation
22 methods arises naturally.

23 Likelihood-free methods aim at bypassing the above problems by resorting to a simulation-based approach to parameter
24 estimation and model selection (3, 12). A subset of likelihood-free methods includes approximate Bayesian computation (ABC)
25 methods, which aim at preserving the advantages of Bayesian data analysis even when the likelihood function is intractable or
26 practically impossible to compute (10, 12, 13). ABC methods approximate the likelihood function by repeatedly sampling
27 parameters from a pre-specified prior distribution $p(\boldsymbol{\theta})$ and then simulating multiple datasets by running the generative model
28 $q(\boldsymbol{\theta})$ using the sampled parameters. Thus, the core ingredients of ABC methods are a prior on $\boldsymbol{\theta}$, and a generative model $q(\boldsymbol{\theta})$,
29 usually specified as a function code in a general-purpose programming language (10, 14).

30 Performing approximate inference comes at the cost of incurring additional approximation error, which accumulates on
31 top of the irreducible estimation error. Within the context of approximate inference, the most common manifestations of
32 approximation error include: *i*) imprecise form of the posterior; *ii*) imprecise posterior moments; *iii*) under- or overestimation
33 of uncertainty. Different approximation methods usually involve multiple trade-offs between minimizing approximation error
34 and keeping computational time within reasonable bounds (3, 15).

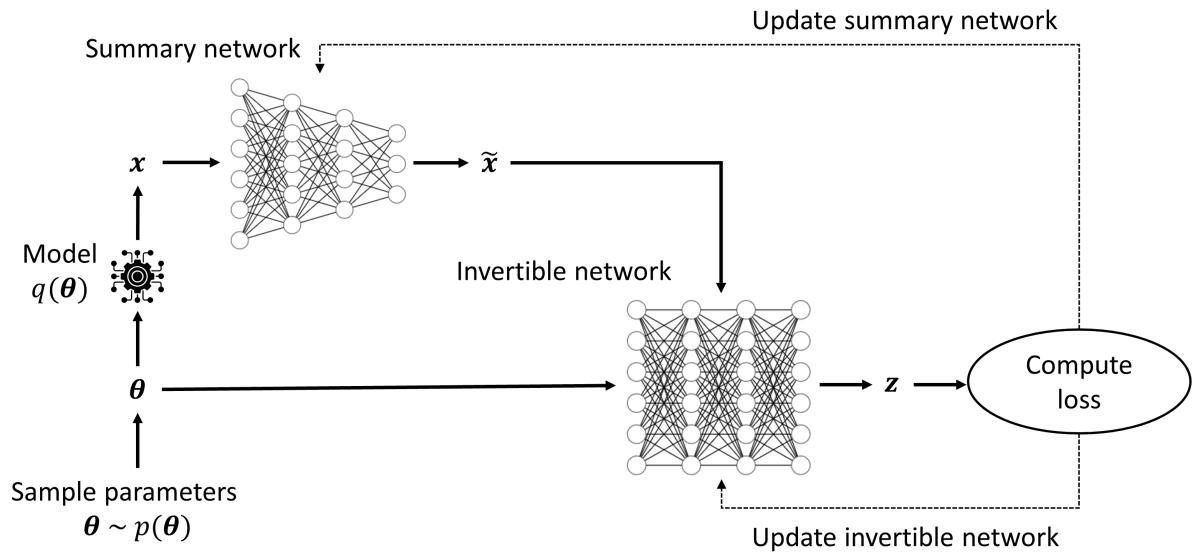
35 Recently, ideas from machine learning and deep learning research have entered the field of likelihood-free inference in
36 an attempt to overcome some of the shortcomings of traditional methods (5, 6, 16–20). The most common approach has
37 been to cast the problem of parameter estimation as a supervised learning task. In this setting, a large dataset of the form
38 $\mathcal{D} = \{h(\mathbf{x}^{(i)}), \boldsymbol{\theta}^{(i)}\}_{i=1}^n$ is created by repeatedly sampling from $p(\boldsymbol{\theta})$ and simulating an artificial dataset \mathbf{x} by running $q(\boldsymbol{\theta})$ with
39 the sampled parameters. Usually, the dimensionality of the simulated data is reduced by computing summary statistics with
40 a fixed summary function $h(\mathbf{x})$. Then, a supervised learning algorithm $f(h(\mathbf{x}); \boldsymbol{\phi}) = \hat{\boldsymbol{\theta}}$ with learnable parameters $\boldsymbol{\phi}$ (e.g.,
41 linear regression, random forest, neural network) is trained on the simulated data to later output an estimate of the true data
42 generating parameters. Thus, $f(h(\mathbf{x}); \boldsymbol{\phi})$ essentially attempts to “learn” the intractable inverse model $\boldsymbol{\theta} = q^{-1}(\mathbf{x})$.

43 Inspired by previous machine learning approaches, the current work proposes a novel and universal likelihood-free method
44 capable of performing full Bayesian inference on any mathematical process model from which simulations can be obtained. It
45 treats parameter inference as a task of inverting a generative model and achieves this by drawing on the modern framework
46 of deep probabilistic modeling for tackling intractable posteriors (21–24). The method integrates two separate deep neural
47 networks modules (detailed in the **Methods** section; see Figure 1) trained jointly on simulated data: a *summary network* and
48 an *invertible network*.

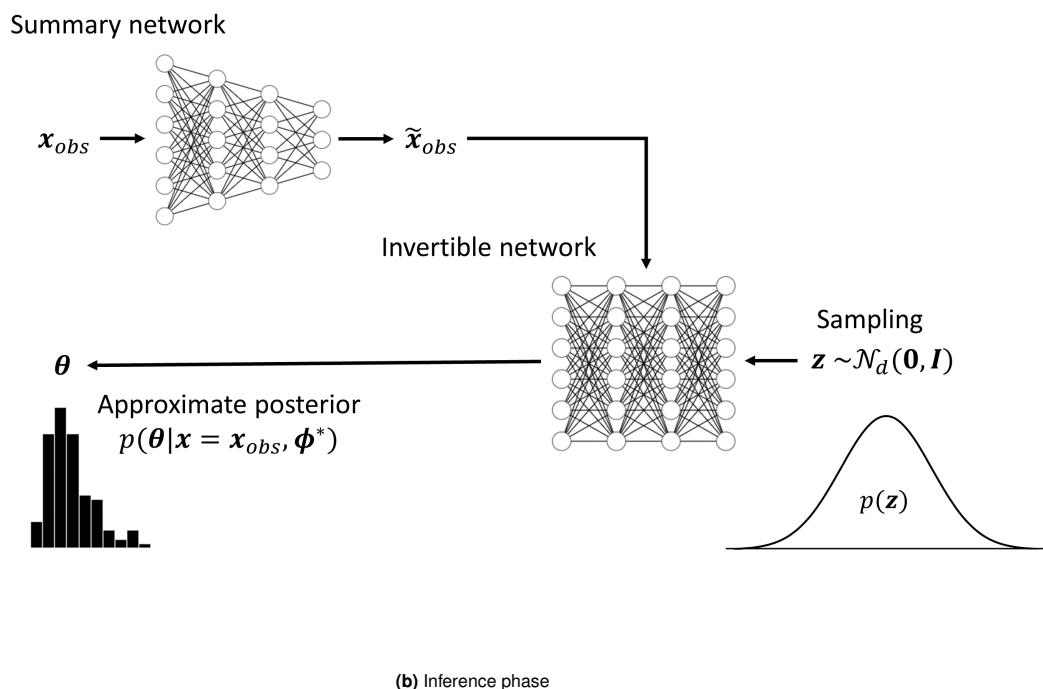
49 The *summary network* is responsible for learning the most informative summary statistics directly from data. It should be
50 designed to follow the functional and probabilistic symmetries inherent in the data, e.g. a permutationally invariant network for
51 *i.i.d.* data (25), a recurrent network for time-series data, or a convolutional network for grid-like data (26). The computation
52 of summary statistics is a crucial aspect in likelihood-free inference. Previous approaches mainly use hand-crafted summary
53 statistics tailored to the specific application. However, in many application, it is not straightforward to settle upon a set of
54 good summary statistics. Thus, the summary network completely eliminates the need to manually specify a fixed number of
55 summary statistics and makes the method independent of the format or the size of the data.

56 The *invertible network* is responsible for learning the posterior of the model parameters given the observed and summarized
57 data. It is based on the recently developed flow-based architecture (22–24). Flow-based methods provide exact latent-variable
58 inference and log-likelihood evaluation when operating at optimum. In the **Methods** section, we show that our method
59 maximizes the posterior over model parameters directly when cast in the context of likelihood-free inference. Furthermore,
60 flow-based methods are capable of approximating very high-dimensional distributions (e.g., the pixels of an image). Once
61 trained with a sufficient amount of simulated data, our invertible network can perform full Bayesian inference on large amounts
62 of real data from a given domain in a single pass. The joint training of a summary network and an invertible network results
63 in a powerful and universal parameter estimation machine capable of inverting complicated statistical problems in various
64 scientific domains (Figure 1). Moreover, the method addresses many of the limitations of previous likelihood-free methods.
65 First, it involves no costly MCMC or rejection sampling, which makes the inference phase lightning fast, as we also make use of
66 GPU-accelerated computation. Second, it involves no fixed summary statistics or kernels, but learns the most informative
67 representation of the data in an end-to-end manner. Third, the method is fully Bayesian, as it directly learns the posterior over
68 model parameters and thus allows for the quantification of uncertainty, which is a crucial requirement in parameter estimation
69 (27, 28). Last, the trained networks can be shared and reused by multiple researches within a scientific domain, thus removing
70 the need for wasteful computations and fitting a separate model for each and every dataset. This pooling of computational
71 resources across researches is an important step forward in mathematical modeling, as has been recently argued (17).

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



(a) Training phase



(b) Inference phase

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.

80 (29) for creating and training custom invertible networks with GPU support, along with some validation tools.

81 Methods

82 **Notation.** In the following, we denote observed or simulated univariate datasets from the mathematical model of interest as
83 $\mathbf{x} = (x_1, x_2, \dots, x_n)$, and multivariate datasets as $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$. The parameters of a mathematical model are represented
84 as a vector $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_d)$, and all trainable parameters of the invertible and summary neural networks as $\phi = (\phi_{inv}, \phi_{sum})$.
85 The number of parameters of a mathematical model will be denoted as d , and the number of simulated data points as n .

86 **Deep Probabilistic Modeling.** Our method draws on major advances in modern deep probabilistic modeling, also referred to
87 as deep generative modeling (21, 22, 25, 30). A hallmark idea in deep probabilistic modeling is to handle intractable target
88 probability distributions by sampling from simpler distributions (e.g., Gaussian or uniform distributions) and transforming these
89 samples via a complex non-linear, learnable transformations. Most popular deep probabilistic models entail two phases. During
90 the *training phase*, a transformation from the simple to the desired target distribution is learned by optimizing a cost function
91 via backpropagation (see Figure 1a). During the *inference phase*, samples from the target distribution are obtained by sampling
92 from the simple distribution and applying the transformation learned during the training phase (see Figure 1b). Using this
93 approach, recent applications of deep probabilistic models have achieved unprecedented results on extremely high-dimensional
94 and intractable problems (e.g., complex data distributions such as natural images, music, or text).

95 In the context of mathematical modeling and Bayesian inference, the target distribution is the posterior distribution of model
96 parameters $p(\boldsymbol{\theta}|\mathbf{x})$ capturing our uncertainty about the numerical values of parameters given empirical data. We can leverage
97 the fact that most mathematical models are generative in nature and as such can be used to perform multiple simulations
98 of the process of interest. By specifying a prior distribution over the model parameters $p(\boldsymbol{\theta})$, one can generate arbitrarily
99 large datasets of the form $\mathbf{D} = \{\mathbf{x}^{(i)}, \boldsymbol{\theta}^{(i)}\}_{i=1}^n$ and use a deep generative model to learn a probabilistic mapping from data
100 to parameters. Thus, at inference time, one can condition the model on observed data \mathbf{x}_{obs} and obtain samples $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_L$
101 approximating the posterior $p(\boldsymbol{\theta}|\mathbf{x} = \mathbf{x}_{obs})$ in the manner described above.

102 In the current work, we propose to implement and use a conditional invertible neural network (cINN) architecture. Previously,
103 INNs have been successfully employed to model data from astrophysics and medicine (21). We adapt the model to suit the task
104 of parameter estimation in the context of mathematical modeling (see Figure 1 for a full graphical illustration of the method)
105 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) (21, 22, 24).
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:

$$v_1 = u_1 \odot \exp(s_1(u_2)) + t_1(u_2) \quad [2]$$

$$v_2 = u_2 \odot \exp(s_2(v_1)) + t_2(v_1) \quad [3]$$

The outputs $\mathbf{v} = (v_1, v_2)$ are then concatenated again and passed to the next ACB. The inverse operation is given by:

$$u_2 = (v_2 - t_2(v_1)) \odot \exp(-s_2(v_1)) \quad [4]$$

$$u_1 = (v_1 - t_1(u_2)) \odot \exp(-s_1(u_2)) \quad [5]$$

106 An additional property of this design, which becomes relevant later for optimization, is that the operations of the ACB have
107 tractable, and cheaply computable Jacobians (strictly upper or lower triangular matrices). Furthermore, the internal networks
108 $s_1(\cdot), s_2(\cdot), t_1(\cdot), t_2(\cdot)$ can be represented by arbitrarily complex neural networks, which themselves need not be invertible, since
109 they are only ever evaluated in the forward direction during both the forward and the inverse pass through the ACB. To ensure
110 that the model is powerful enough to represent complicated distributions, we chain multiple ACBs, so that the output of each
111 ACB becomes the input of the next (see Figure 2b). In this way, the whole chain remains invertible from the first input to the
112 last output and can be viewed as a single function parameterized by trainable parameters ϕ_{inv} .

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

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

together with the inverse operation:

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

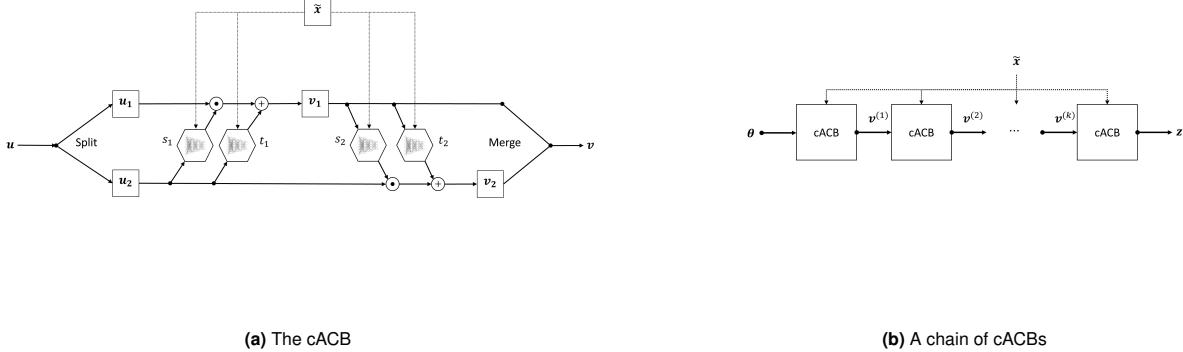


Fig. 2. A diagram of the conditional version of the affine coupling block (cACB). **(a)** Each cACB consists of four internal networks performing the invertible operations described in the text; **(b)** In practice, we chain multiple cACBs to obtain higher representational capacity. Each cACB layer uses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

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

117 **Summary Network.** Since in practice the conditioning data set \mathbf{x} can have variable number of input points (e.g., trial sizes,
 118 time points) and exhibit various redundancies, the cINN can profit from some form of dimensionality reduction applied to the
 119 data. Ideally, we want to avoid hand-crafted summary statistics, and instead learn the most informative summary statistics
 120 directly from data. Therefore, instead of feeding the raw simulated (observed) data to each ACB, we pass the data through
 121 an additional summary network to obtain a fixed-sized vector of learned summary statistics $\tilde{\mathbf{x}} = h(\mathbf{x}; \phi_{sum})$ and learn the
 122 parameters of the summary network h jointly with those of the cINN chain via backpropagation. Thus, the current method
 123 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:

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

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

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

Note, that by maximizing Eq. 9 we are maximizing the posterior over model parameters of interest $p(\theta|\mathbf{x}, \phi)$. 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 θ to that of \mathbf{z} via the change of variable theorem:

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

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

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 log 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 :

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

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

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

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

126 Minimizing Eq.14 can be interpreted as searching for the optimal neural network parameters ϕ^* which maximize the
 127 probability of model parameters θ given data x . This is exactly the probability we are concerned with in Bayesian inference.
 128 Note that our formulation maximizes the posterior of model parameters directly, in contrast to variational methods which
 129 optimize a lower bound on the posterior (20, 30). Once the backpropagation algorithm has settled to a local minimum of the ML
 130 loss, one can easily obtain samples from the approximate posterior $p(\theta|x = x_{obs}, \phi = \phi^*)$, based on an observed dataset x_{obs} , by
 131 repeatedly sampling $z^{(l)} \sim \mathcal{N}_d(\mathbf{0}, I)$, and then passing $z^{(l)}$ in reverse to the cINN in order to compute $\theta^{(l)} = f^{-1}(z^{(l)}; x_{obs}, \phi^*)$
 132 for $l = 1, \dots, L$. Figure 1b illustrates the inference phase of the method. It is worth noting that sampling a large number of
 133 parameter values from the approximate posterior takes a negligible amount of time, since it only requires a single pass through
 134 the cINN in reverse.

135 **Training the Networks.** We train all cINNs and summary networks described in this paper jointly via backpropagation. For all
 136 following experiments, we use the Adam optimizer with a starter learning rate of 10^{-3} and an exponential decay rate of .95.
 137 We set the weight regularization parameter τ to a value of 10^{-5} . We perform 50000 to 100000 iterations (network updates) for
 138 each of the examples in this paper, and report the results on the trained network. Note, that we did not perform an extensive
 139 search for optimal values of the cINN hyperparameter, but use a cINN with 10 ACBs all examples in this paper (see SI for
 140 details of the cINN). All networks were trained on a single-GPU machine equipped with NVIDIA GTX1060 graphics card.
 141 Regarding the data generation step, we use two training approaches.

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

148 The second approach follows a different strategy, which is used mainly in the field of active learning. Correspondingly,
 149 a dataset, or a batch of datasets, is created on the fly and then passed through the neural network. This training regime
 150 has the advantage that the network never “experiences” the same input data twice. Moreover, training can continue as long
 151 as the network keeps improving (i.e., the loss keeps decreasing), since overfitting in the classical sense is nearly impossible.
 152 However, if the simulations are computationally expensive and researchers need to experiment with different networks or
 153 training hyperparameters, it might be beneficial to switch to the first regime, since simulation and training in the active
 154 learning regime are tightly intertwined.

155 **Putting It All Together.** On an abstract level, our method requires three key ingredients: 1) a mathematical process model
 156 $q(\theta)$ capable of simulating data x ; 2) a prior distribution over the model parameters $p(\theta)$ encoding our prior beliefs about
 157 plausible parameter values; 3) an invertible neural network f_ϕ capable of approximating a large enough family of probability
 158 distributions (see Figure 2). In practice, a chain of up to 10 ACBs should suffice to learn most distributions encountered in the
 159 life sciences, since they tend to be unimodal and relatively simple (in contrast to the distribution of natural images or words in
 160 a spoken language). From these three ingredients, a universal and reusable sampler can be designed for likelihood-free Bayesian
 161 estimation of both tractable and intractable mathematical models. **Algorithm 1** describes the essential steps of the method
 162 using an arbitrary summary network and employing the active learning training regime.

Algorithm 1 Bayesian likelihood-free inference with invertible neural networks

```

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

163 The backpropagation algorithm works by computing the gradients of the loss function w.r.t. the parameters of the neural

164 networks and then adjusting the parameters, so as to drive the loss function to a local minimum. We experienced no instability
 165 or convergence issues during training with the ML loss. Note, that steps 12 – 15 of **Algorithm 1** can be executed in parallel
 166 with GPU support.

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

171 Results

172 We consider a number of metrics for determining the performance of our method. To assess the goodness of parameter
 173 recovery, we compute the normalized root mean squared error (NRMSE) and the coefficient of determination (R^2) metrics of
 174 the parameter estimates. To assess the recovery of the full posterior, we compute the Kullback-Leibler (KL) divergence (31)
 175 between the true and the approximate distributions for the toy example, and use simulation-based calibration (SBC, (32)) for
 176 the other examples where the analytic posterior is not available in closed-form. Details for computing the *NRMSE*, R^2 , and
 177 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 setting, 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 model each $y^{(i)}$ as being conditionally Gaussian given $x^{(i)}$, i.e., $y^{(i)} \sim \mathcal{N}(\boldsymbol{\theta}^T \mathbf{x}^{(i)}, a^{-1})$ 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:

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

where \mathbf{X} denotes the design matrix containing all $\mathbf{x}^{(i)}$ stacked row-wise. 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 $p(\boldsymbol{\theta}|\mathbf{D}) = \mathcal{N}_d(\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$ where $\boldsymbol{\Lambda}$ denotes the posterior precision matrix (inverse covariance matrix), $\boldsymbol{\mu}$ the posterior mean vector, which are computed as follows:

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

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

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

180 For the following application, we set $d = 4$, and $a = b = 1$. The design matrices for each iteration contain a variable number
 181 n of *i.i.d.* data points drawn from a unit Gaussian $\mathbf{x}^{(i)} \sim \mathcal{N}_d(\mathbf{0}, \mathbf{I})$ for $i = 1, \dots, n$. The number of trials is drawn from a
 182 uniform distribution $n \sim \mathcal{U}(50, 500)$ at each training iteration (Lines 2–9 of **Algorithm 1**).

183 The results on the toy Bayesian regression are depicted in **Figure 3**. The approximate posterior means show negligible
 184 deviations from the analytic posterior means as quantified by very small NRMSEs (as small as 0.002) and very high R^2 (as
 185 high as 1.0) over all n sampled during training. This suggests near-perfect estimation of the true posterior means. Further, the
 186 estimates become increasingly more accurate, as the number of data points n increases. An inspection of the posterior variances
 187 over the different n reveals that the estimated variance follows closely the decrease in analytic variance with increasing n .
 188 However, the analytic variance is slightly underestimated at smaller n and slightly overestimated at larger n . This pattern is
 189 also revealed by the KL divergence plot (see **SI**). This result might be attributable to an underexpressive summary network;
 190 another possibility is that the networks need to be trained longer with smaller learning rate decay.

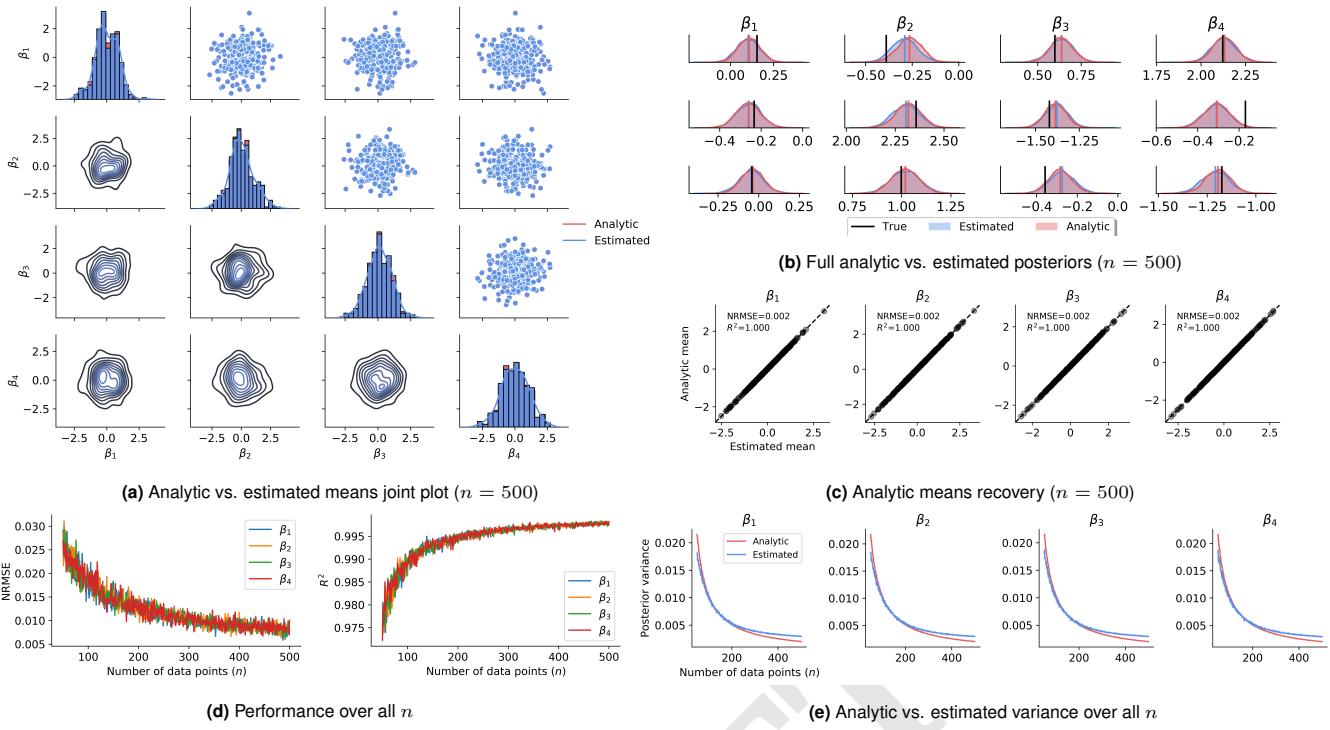


Fig. 3. Results on the Bayesian toy regression example. (a) Analytic vs. estimated posterior means on the validation sample. We observe a near-perfect overlap between analytic and estimated posterior means and no spurious covariances; (b) Some example draws from the estimated posteriors. Visual inspections 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:

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

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

for $t = 1, \dots, T$ where N_t is the expected number of individuals at time t , r is the growth rate, ρ is a scaling parameter and $\epsilon_t \sim \mathcal{N}(0, \sigma^2)$ is random Gaussian noise (17). The Ricker model has no known likelihood function and is thus 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).

What if the data does not contain information about a given parameter? In this case, any good estimation method should detect this, and return the prior of the particular parameter. To test this, we add a random uniform variable $u \sim \mathcal{U}(0, 1)$ to the parameter vector θ and train the model with this additional dummy variable. Then, we inspect the posterior of u for uniformity.

The results on the Ricker model are depicted in Figure 3. As evident from the graphics, the 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 σ seems to be 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 (32), the approximate posteriors of r and ρ slightly underestimate the true posterior means, whereas the approximate posterior of σ tends to overestimate the true posterior variance. These deviations appear to be due to the fact that recovery worsens at extreme values of the parameters. This is unsurprising, as the data generated with these parameters is highly implausible, which in some cases might even render a model unidentifiable.

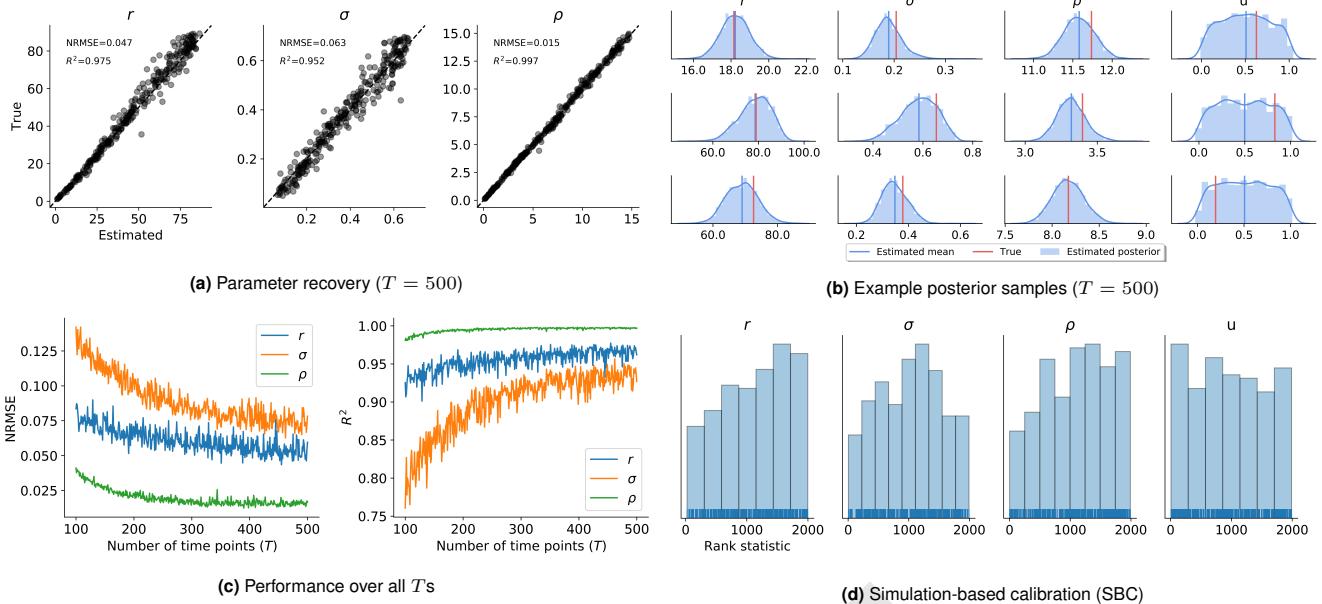


Fig. 4. Results on the Ricker model. **(a)** Parameter recovery for the maximum number of generations used during training ($T = 500$); **(b)** Example posteriors for three test datasets. We observe that the posterior of the uniform noise variable u is equal to the prior, i.e., the method detects that 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 ρ slightly underestimate the true posterior means, whereas the approximate posterior of σ tends to overestimate the true posterior variance.

Example 2 - The Lévy-Flight Model. Evidence accumulator models (EAMs) describe (perceptual) decision making by a set of neurocognitively motivated parameters (33). 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 and 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 exceeds a threshold. In its most general formulation, the forward model of EAMs takes the form of a stochastic differential equation given by (4):

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

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

EAMs are particularly amenable for likelihood-free inference, since most members of the family turn out to be intractable (34). 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, (35)). 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 from the literature. The forward equation of the LFM is given by:

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

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

The LFM has three additional parameters: a 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 RT data from two experimental conditions with two different drift rates (see SI for details of the simulation). 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.

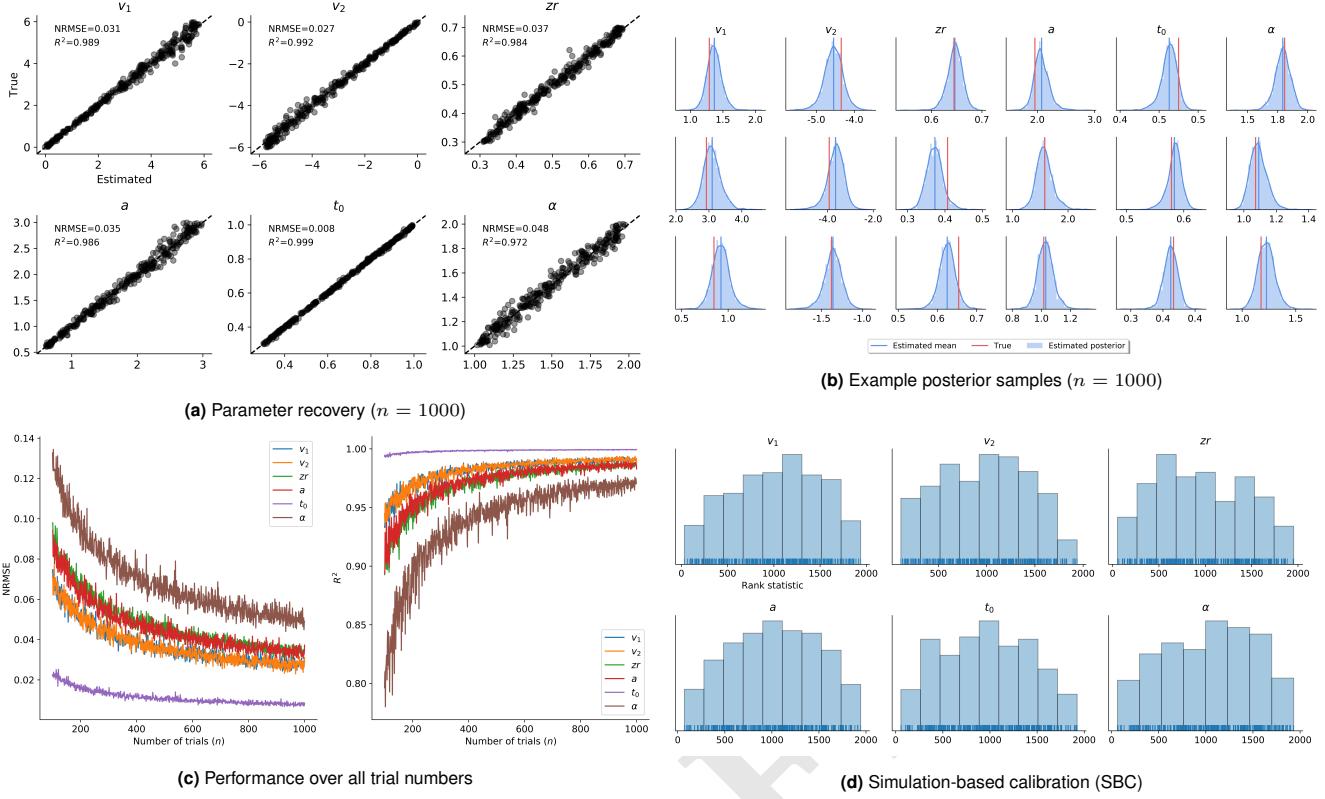


Fig. 5. Results on the LFM model. **(a)** Parameter recovery for the maximum number of trials used during training ($n = 1000$); **(b)** Example posteriors for three test datasets; **(c)** NRMSE and R^2 performance metrics over all n trials used during training. Again, we observe that the recovery remains good overall, 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 4. To our knowledge, this is the first Bayesian treatment of the LFM, as its intractability makes traditional methods prohibitively slow. We observe state-of-the art 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 α 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 (32).

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 (8, 9, 36). 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:

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

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

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

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

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

where $S + I + R = N$ give the number of susceptible, infected, and recovered individuals, respectively. The parameter β controls the transition rate from being susceptible to infected, and γ controls the transition rate from being infected to recovered. The number of individuals moving from S to I , given by ΔN_{SI} , and the number of people moving from I to R , given by ΔN_{IR} , over a time interval Δt are modeled as binomial random variables. The above listed stochastic system has no known analytic solution and thus requires numerical simulation methods for finding optimal parameter values. Cast as a parameter estimation task, the challenge is to recover $\theta = \{\beta, \gamma\}$ from three dimensional time-series data $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T)$ 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 .

237 During training of the networks, we simulate time-series from the stochastic SIR model with varying lengths. The number
238 of time points T is drawn from a uniform distribution $T \sim \mathcal{U}(100, 200)$ at each training iteration (see the **SI** for more details
239 about the simulation). This range is chosen since we observed that the system reaches an equilibrium in almost all parameter
240 settings after $T > 200$.

241 The results on the SIR model are depicted in Figure XXX.

242 **Example 4 - Single-Cell RNA Sequencing.** Single-cell RNA sequencing (scRNA-seq) is a method to uncover the dynamics of
243 gene expression within single cells (1, 37). Simulation models for scRNA-seq attempt to mimic the process of sequencing
244 real data samples by combining statistical and algorithmic procedures. The recently developed *Splat* simulation model (1)
245 implements a hierarchical model where the mean expression levels for multiple genes are samples from a Gamma distribution,
246 and the number of times a gene has been sequenced in a cell is then sampled from a Poisson distribution. The output of the
247 *Splat* simulation is thus a matrix \mathbf{X} of *Gene x Cell* counts. Figure XXX depicts the *Splat* algorithm for simulating scRNA-seq
248 data.

249 As a last example for our method, we use simulations from the *Splat* model and attempt to recover the data-generating
250 parameters. We simulate 250000 *Gene x Cell* count matrices with different parameter settings (see **SI** for details of the
251 simulation) and discard implausible simulations (e.g., matrices where more than 90% of the counts are 0). We train the
252 networks for 100 epochs through the entire dataset and evaluate the performance on a separate validation set or matrices.

253 Figure XXX depicts the results...

254 Discussion

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

261 Our method combines the universal approximation power of deep learning methods with the important uncertainty
262 quantification assets of Bayesian inference (27, 28). Besides being capable of performing rapid Bayesian inference on intractable
263 mathematical models, our method provides a general framework for designing reusable “parameter estimation machines” for
264 various research domains. Moreover, the method is not confined solely to inference on intractable models, but can also prove as
265 a viable alternative in modeling contexts where standard Bayesian or frequentist inference methods are available, but inference
266 is nevertheless prohibitively slow.

267 Inspired by previous machine learning approaches to likelihood-free inference (5, 6, 16–20), our method shares many of the
268 advantages of these methods and further overcomes some important limitations.

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

275 Second, we showed that the ML loss optimizes directly the posterior over model parameters of interest, which is in contrast
276 to ELBO-based methods which optimize a lower-bound on the posterior (20, 30). Thus, inference is exact when the ML
277 is minimized (22, 24). This *optimal performance* claim is confirmed by our toy Bayesian regression example, in which we
278 observe negligibly small deviations of the approximate from the true posterior. Further, researchers are often interested in
279 some summary of the posterior, for instance, the posterior mean or the posterior variance (16, 18). We demonstrated that
280 our method exhibits excellent recovery of the posterior means throughout the examples. We also showed that the recovery
281 becomes better with increasing number of observed data points, while the variance of the estimates becomes larger. This is
282 an important and highly desirable property of any parameter estimation method, as it mirrors the increase in information
283 following increasing number of data points.

284 Third, the largest computational cost of our method is paid during training. Once trained, the networks can be used and
285 reused to perform inference on large numbers of datasets within seconds and across a given research domain. Indeed, there
286 are many instances of research domains where a single model is extensively explored and independently fitted by multiple
287 researchers to test scientific hypotheses (1, 33, 35?). These research domains are expected to benefit the most from learning the
288 “model universe” once and then inverting the model multiple times for rapid inference on different datasets. In this regard, our
289 method is similar to the recently introduced prepaid method (17) which uses a database of pre-computed summary statistics
290 and nearest-neighbors for inference. Note, however, that our method does not need to store training data on disk, since the
291 “knowledge” about the relationship between data and parameters is compressed into the networks’ weights. Moreover, all
292 computations involved in our method benefit from a high degree of parallelism and can thus utilize the advantages of modern
293 GPUs.

294 These advantages notwithstanding, some limitations of the method deserve a mention. Even though high-level deep learning
295 libraries, such as *TensorFlow* or *Torch*, allow for rapid and relatively straightforward development of various neural network

296 architectures, the implementational burden associated with the current method is still reasonably high. In order to ease the
297 understanding and independent application of the method, we provide fully functioning code to reproduce and study all of the
298 examples tackled in this paper (<https://github.com/stefanradev93/cINN>). In addition, we also provide implementation of all
299 tools for performance validation used throughout the paper. Moreover, we are currently developing a general user-friendly
300 software, which should abstract away most of the methodological complexities from the user. Another potential shortcoming of
301 the method is the seemingly overwhelming number of hyperparameters that might require fine-tuning by the user for optimal
302 performance on a given task. However, we observe that many of the default hyperparameter values are sufficient to achieve
303 excellent performance, and starting with a relatively large default network of 10 ACBs does not appear to hurt performance or
304 destabilize training, even if the model to be learned is relatively simple. We expect that a single architecture should be able to
305 perform well on almost all models from a given domain (i.e., a single architecture for decision-making models). Future research
306 should investigate the question of generality by applying the method to challenging parameter estimation tasks across different
307 research domains.

308 Conclusion

309 As formal theories in various scientific disciplines (especially in the younger sciences, such as, neuroscience, cognitive science,
310 computational biology, etc.) become increasingly complex, the need for powerful and universally applicable likelihood-free
311 estimation methods becomes increasingly pressing. In the present work, we addressed this need by introducing a method
312 potentially applicable to *any* modeling scenario in *any* research domain where simulations from the process model can be
313 obtained. We hope that the new method will enable researchers from a variety of fields to accelerate model-based inference and
314 will further prove its utility beyond the examples considered in this paper.

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

- 317 1. Zappia L, Phipson B, Oshlack A (2017) Splatter: simulation of single-cell rna sequencing data. *Genome biology* 18(1):174.
- 318 2. Beaumont MA, Zhang W, Balding DJ (2002) Approximate bayesian computation in population genetics. *Genetics* 162(4):2025–2035.
- 319 3. Palestro JJ, Sederberg PB, Osth AF, Van Zandt T, Turner BM (2018) *Likelihood-free methods for cognitive science*. (Springer).
- 320 4. Usher M, McClelland JL (2001) The time course of perceptual choice: the leaky, competing accumulator model. *Psychological review* 108(3):550.
- 321 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*.
- 322 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.
- 323 7. Wood SN (2010) Statistical inference for noisy nonlinear ecological dynamic systems. *Nature* 466(7310):1102.
- 324 8. Keeling MJ, Rohani P (2011) *Modeling infectious diseases in humans and animals*. (Princeton University Press).
- 325 9. Hethcote HW (2000) The mathematics of infectious diseases. *SIAM review* 42(4):599–653.
- 326 10. Csillér K, Blum MG, Gaggiotti OE, François O (2010) Approximate bayesian computation (abc) in practice. *Trends in ecology & evolution* 25(7):410–418.
- 327 11. Toni T, Stumpf MP (2009) Simulation-based model selection for dynamical systems in systems and population biology. *Bioinformatics* 26(1):104–110.
- 328 12. Turner BM, Sederberg PB (2014) A generalized, likelihood-free method for posterior estimation. *Psychonomic bulletin & review* 21(2):227–250.
- 329 13. Sunnåker M, et al. (2013) Approximate bayesian computation. *PLoS computational biology* 9(1):e1002803.
- 330 14. 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.
- 331 15. Frazier DT, Martin GM, Robert CP, Rousseau J (2018) Asymptotic properties of approximate bayesian computation. *Biometrika* 105(3):593–607.
- 332 16. 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*.
- 333 17. Mestdagh M, Verdonck S, Meers K, Loosens T, Tuerlinckx F (2018) Prepaid parameter estimation without likelihoods. *arXiv preprint arXiv:1812.09799*.
- 334 18. Raynal L, et al. (2018) Abc random forests for bayesian parameter inference. *Bioinformatics* 35(10):1720–1728.
- 335 19. 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.
- 336 20. Papamakarios G, Murray I (2016) Fast ϵ -free inference of simulation models with bayesian conditional density estimation in *Advances in Neural Information Processing Systems*. pp. 1028–1036.
- 337 21. Ardizzone L, et al. (2018) Analyzing inverse problems with invertible neural networks. *arXiv preprint arXiv:1808.04730*.
- 338 22. Kingma DP, Dhariwal P (2018) Glow: Generative flow with invertible 1x1 convolutions in *Advances in Neural Information Processing Systems*. pp. 10215–10224.
- 339 23. 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*.
- 340 24. Dinh L, Sohl-Dickstein J, Bengio S (2016) Density estimation using real nvp. *arXiv preprint arXiv:1605.08803*.
- 341 25. Bloem-Reddy B, Teh YW (2019) Probabilistic symmetry and invariant neural networks. *arXiv preprint arXiv:1901.06082*.
- 342 26. Goodfellow I, Bengio Y, Courville A (2016) *Deep learning*. (MIT press).
- 343 27. 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.
- 344 28. Gelman A, et al. (2013) *Bayesian data analysis*. (Chapman and Hall/CRC).
- 345 29. 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.
- 346 30. Kingma DP, Welling M (2014) Auto-encoding variational bayes. *stat* 1050:1.
- 347 31. Hershey JR, Olsen PA (2007) Approximating the kulback 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.
- 348 32. Talts S, Betancourt M, Simpson D, Vehtari A, Gelman A (2018) Validating bayesian inference algorithms with simulation-based calibration. *arXiv preprint arXiv:1804.06788*.
- 349 33. Ratcliff R, McKoon G (2008) The diffusion decision model: theory and data for two-choice decision tasks. *Neural computation* 20(4):873–922.
- 350 34. 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.
- 351 35. 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.
- 352 36. 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.
- 353 37. Ozsolak F, Milos PM (2011) Rna sequencing: advances, challenges and opportunities. *Nature reviews genetics* 12(2):87.