

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 20, 2019

Parametric models of complex processes are ubiquitous throughout the sciences. As the processes under study and the models describing them become increasingly complex, parameter estimation with standard Bayesian and frequentist methods can quickly become intractable. To address this, we propose a novel method for likelihood-free inference based on invertible neural networks. The method is capable of performing fast full Bayesian inference on large amounts of data by training the networks on simulated data and learning to invert the model under study. The method is independent of particular data formats, as it includes a summary network trained to embed the observed data into fixed-size vectors in a data-driven way. This makes the method applicable to various scenarios where standard inference techniques fail. We demonstrate the utility of the method on a toy model with known analytic posterior and on example models from population dynamics, epidemiology, cognitive science and genetics. We argue for a general framework for building reusable parameter estimation machines for 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 their most abstract form, mathematical models are specified by a set of latent parameters θ and a generative model q mapping parameters to manifest quantities x :

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

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

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

Moreover, as mathematical models and processes under description become increasingly complex, parameter estimation and model selection can quickly become intractable with standard Bayesian and frequentist method. Complex models specified by a generative stochastic mechanism do not always provide a closed-form solution for the *likelihood function* (3, 9, 10). This poses great difficulties for Bayesian and frequentist methods alike, since both depend explicitly on the numerical evaluation of

Significance Statement

Authors must submit a 120-word maximum statement about the significance of their research paper written at a level understandable to an undergraduate educated scientist outside their field of speciality. The primary goal of the Significance Statement is to explain the relevance of the work in broad context to a broad readership. The Significance Statement appears in the paper itself and is required for all research papers.

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.radevpsychologie.uni-heidelberg.de

a likelihood function as a proxy for assessing model fit to data. Even if a likelihood function is available, inference could be prohibitively slow for real-world applications. Therefore, the need for powerful and reliable likelihood-free estimation methods arises naturally.

Likelihood-free methods aim at bypassing the intractability problem by resorting to a simulation-based approach to inference and model selection (3). 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, 12). ABC methods approximate the likelihood function by repeatedly sampling parameters from a pre-specified prior distribution $\pi(\theta)$ and then simulating multiple datasets by running the generative model $q(\theta)$ using the sampled parameters. Thus, the core ingredients of ABC methods are a prior on θ , and a generative model $q(\theta)$, usually specified as a function code in a general-purpose programming language (9, 13).

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, 14).

Recently, ideas from machine learning and deep learning research have entered the field of likelihood-free inference in an attempt to overcome some of the shortcomings of traditional methods (5, 6, 15–19). 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 $D = \{h(\mathbf{x}^{(i)}), \theta^{(i)}\}_{i=1}^n$ is created by repeatedly sampling from $\pi(\theta)$ and simulating an artificial dataset \mathbf{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 $h(\mathbf{x})$. Then, a supervised learning algorithm $f(h(\mathbf{x}); \phi) = \hat{\theta}$ with learnable parameters ϕ (e.g., linear regression, random forest, neural network) is trained on the simulated data to later output an estimate of the true data generating parameters. Thus, $f(h(\mathbf{x}); \phi)$ essentially attempts to “learn” the intractable inverse model $\theta = q^{-1}(\mathbf{x})$.

Inspired by previous machine learning approaches, 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 a generative model and achieves this by drawing on the modern framework of deep probabilistic modeling for tackling intractable posteriors (20–23). The method integrates two separate deep neural networks modules (detailed in the **Methods** section) 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 (24), a recurrent network for time-series data, or a convolutional network for grid-like data (25). The computation of summary statistics is a crucial aspect in likelihood-free inference. Previous approaches mainly use hand-crafted summary statistics tailored to the specific application. However, in many application, it is not straightforward to settle upon a set of good summary statistics. Thus, the summary network completely eliminates the need to manually specify a fixed number of summary statistics and 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 (21–23). 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 very high-dimensional distributions (e.g., images). Once trained with a sufficient amount of simulated data, our invertible network can perform full Bayesian inference on large amounts of real data from a given domain in a single pass.

The joint training of a summary network and an invertible network results in a powerful and universal parameter estimation machine capable of inverting complicated statistical problems in various scientific domains (see Fig.1). Moreover, the method addresses many of the limitations of previous likelihood-free methods. First, it involves no costly MCMC or rejection sampling, which makes the inference phase lightning fast, as we also make use of GPU-accelerated computation. Second, it involves no fixed summary statistics or kernels, but learns the most informative representation of the data in an end-to-end manner. Third, the method is fully Bayesian, as it directly learns the posterior over model parameters and thus allows for the quantification of uncertainty, which is a crucial requirement in parameter estimation (26, 27). Last, the trained networks can be shared and reused by multiple researches within a scientific domain, thus removing the need for wasteful computations and fitting a separate model for each and every dataset. This pooling of computational resources across researches is an important step forward in mathematical modeling, as has been recently argued (16).

To illustrate the utility of the new method, we first apply it to a toy Bayesian regression model with known posterior. Then, we present applications to intractable models from cognitive science, population dynamics, epidemiology, and genetics and demonstrate state-of-the-art parameter recovery. The outline of the remaining manuscript is as follows: The **Methods** section introduces the main building blocks of the new method and summarizes the main steps as pseudocode. The **Results** section presents the various applications of the model to real-world research domains. Finally, the **Discussion** section lists the advantages of the current method, treats some potential pitfalls and explores future research vistas. Python code and data for all applications are freely available as Jupyter notebooks at <https://github.com/stefanradev93/cINN> and as a small library based on *TensorFlow* (28) for creating and training custom invertible networks with GPU-support.

Methods

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

Deep Probabilistic Modeling. Our method draws on major advances in modern deep probabilistic modeling, also referred to as deep generative modeling (20, 21, 24, 29). A hallmark idea in deep probabilistic modeling is to handle intractable target probability distributions by sampling from simpler distributions (e.g., Gaussian or uniform distributions) and transforming these samples via a complex non-linear, learnable transformations. Most popular deep probabilistic models entail two phases. During the *training phase*, a transformation from the simple to the desired target distribution is learned by optimizing a cost function via backpropagation. During the *inference phase*, samples from the target distribution are obtained by sampling from the simple distribution and applying the transformation learned during the training phase. Using this approach, recent applications of deep probabilistic models have achieved unprecedented results on extremely high-dimensional and intractable problems (e.g., complex data distributions such as natural images, music, or text).

In the context of mathematical modeling and Bayesian inference, the target distribution is the posterior distribution of model parameters $p(\boldsymbol{\theta}|\mathbf{x})$ capturing our uncertainty about the numerical values of parameters given empirical data. We can leverage the fact that most mathematical models are generative in nature and as such can be used to perform multiple simulations of the process of interest. By specifying a prior distribution over the model parameters $\pi(\boldsymbol{\theta})$, one can generate arbitrarily 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 to parameters. Thus, at inference time, one can condition the model on observed data \mathbf{x}_{obs} and obtain samples $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_L$ approximating the posterior $p(\boldsymbol{\theta}|\mathbf{x} = \mathbf{x}_{obs})$ in the manner described above.

In the current work, we propose to implement and use a conditional invertible neural network (cINN) architecture. Previously, INNs have been successfully employed to model data from astrophysics and medicine(20). We adapt the model to suit the task of parameter estimation in the context of mathematical modeling (see Figure 1 for a full graphical illustration of the method) 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) (20, 21, 23). Each ACB consists of four separate fully connected neural networks denoted as $s_1(\cdot), s_2(\cdot), t_1(\cdot), t_2(\cdot)$. An ACB is specifically designed to be invertible, which means that in addition to a parametric mapping $f_{\boldsymbol{\phi}_{inv}} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is also learns the inverse mapping $f_{\boldsymbol{\phi}_{inv}}^{-1} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ "for free". Denoting the input vector of $f_{\boldsymbol{\phi}_{inv}}$ as \mathbf{u} and the output vector as \mathbf{v} , it follows that $f(\mathbf{u}; \boldsymbol{\phi}_{inv}) = \mathbf{v}$ and $f^{-1}(\mathbf{v}; \boldsymbol{\phi}_{inv}) = \mathbf{u}$. 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:

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

$$\mathbf{v}_2 = \mathbf{u}_2 \odot \exp(s_1(\mathbf{v}_1)) + t_1(\mathbf{v}_1) \quad [3]$$

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:

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

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

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. 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 $\boldsymbol{\phi}_{inv}$.

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:

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

together with the inverse operation:

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

This process can be interpreted as follows: the forward pass maps data-generating parameters to \mathbf{z} -space using conditional information of \mathbf{x} , while the inverse pass maps data points from the \mathbf{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 $f^{-1}(\mathbf{z}; \mathbf{x}, \boldsymbol{\phi}_{inv})$ to the posterior $p(\boldsymbol{\theta}|\mathbf{x})$.

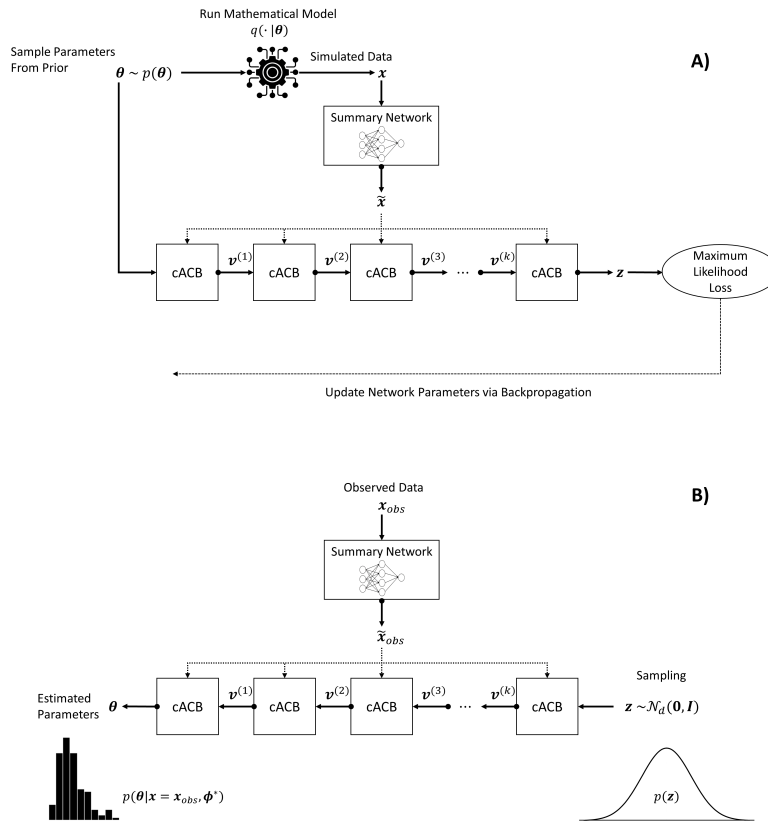


Fig. 1. Placeholder image of a frog with a long example caption to show justification setting.

Summary Network. Since in practice the conditioning data set \mathbf{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 (observed) data to each ACB, we pass the data through an additional summary network to obtain a fixed-sized vector of learned summary statistics $\tilde{\mathbf{x}} = h(\mathbf{x}; \phi_{sum})$ 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.

Submitting Manuscripts. All authors must submit their articles at [PNAScentral](https://www.pnas.org). If you are using Overleaf to write your article, you can use the “Submit to PNAS” option in the top bar of the editor window.

Format. Many authors find it useful to organize their manuscripts with the following order of sections; Title, Author Affiliation, Keywords, Abstract, Significance Statement, Results, Discussion, Materials and methods, Acknowledgments, and References. Other orders and headings are permitted.

Manuscript Length. PNAS generally uses a two-column format averaging 67 characters, including spaces, per line. The maximum length of a Direct Submission research article is six pages and a Direct Submission Plus research article is ten pages including all text, spaces, and the number of characters displaced by figures, tables, and equations. When submitting tables, figures, and/or equations in addition to text, keep the text for your manuscript under 39,000 characters (including spaces) for Direct Submissions and 72,000 characters (including spaces) for Direct Submission Plus.

Data Archival. PNAS must be able to archive the data essential to a published article. Where such archiving is not possible, deposition of data in public databases, such as GenBank, ArrayExpress, Protein Data Bank, Unidata, and others outlined in the Information for Authors, is acceptable.

Language-Editing Services. Prior to submission, authors who believe their manuscripts would benefit from professional editing are encouraged to use a language-editing service (see list at www.pnas.org/site/authors/language-editing.xhtml). PNAS does not take responsibility for or endorse these services, and their use has no bearing on acceptance of a manuscript for publication.

Digital Figures. Only TIFF, EPS, and high-resolution PDF for Mac or PC are allowed for figures that will appear in the main text, and images must be final size. Authors may submit U3D or PRC files for 3D images; these must be accompanied by 2D representations in TIFF, EPS, or high-resolution PDF format. Color images must be in RGB (red, green, blue) mode. Include the font files for any text.

Figures and Tables should be labelled and referenced in the standard way using the `\label{}` and `\ref{}` commands.

Figure 1 shows an example of how to insert a column-wide figure. To insert a figure wider than one column, please use the `\begin{figure*}...\end{figure*}` environment. Figures wider than one column should be sized to 11.4 cm or 17.8 cm wide. Use `\begin{SCfigure*}...\end{SCfigure*}` for a wide figure with side captions.

Single column equations. Authors may use 1- or 2-column equations in their article, according to their preference.

To allow an equation to span both columns, use the `\begin{figure*}...\end{figure*}` environment mentioned above for figures.

Supporting Information Appendix (SI). Authors should submit SI as a single separate PDF file, combining all text, figures, tables, movie legends, and SI references. PNAS will publish SI uncomposed, as the authors have provided it. Additional details can be found here: [policy on SI](#). For SI formatting instructions click [here](#). The PNAS Overleaf SI template can be found [here](#). Refer to the SI Appendix in the manuscript at an appropriate point in the text. Number supporting figures and tables starting with S1, S2, etc.

Authors who place detailed materials and methods in an SI Appendix must provide sufficient detail in the main text methods to enable a reader to follow the logic of the procedures and results and also must reference the SI methods. If a paper is fundamentally a study of a new method or technique, then the methods must be described completely in the main text.

SI Datasets. Supply .xlsx, .csv, .txt, .rtf, or .pdf files. This file type will be published in raw format and will not be edited or composed.

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. Hethcote HW (2000) The mathematics of infectious diseases. *SIAM review* 42(4):599–653.
9. Csilléry K, Blum MG, Gaggiotti OE, François O (2010) Approximate bayesian computation (abc) in practice. *Trends in ecology & evolution* 25(7):410–418.
10. Toni T, Stumpf MP (2009) Simulation-based model selection for dynamical systems in systems and population biology. *Bioinformatics* 26(1):104–110.
11. Turner BM, Sederberg PB (2014) A generalized, likelihood-free method for posterior estimation. *Psychonomic bulletin & review* 21(2):227–250.
12. Sunnåker M, et al. (2013) Approximate bayesian computation. *PLoS computational biology* 9(1):e1002803.
13. 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.
14. Frazier DT, Martin GM, Robert CP, Rousseau J (2018) Asymptotic properties of approximate bayesian computation. *Biometrika* 105(3):593–607.
15. 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*.
16. Mestdagh M, Verdonck S, Meers K, Loossens T, Tuerlinckx F (2018) Prepaid parameter estimation without likelihoods. *arXiv preprint arXiv:1812.09799*.
17. Raynal L, et al. (2018) Abc random forests for bayesian parameter inference. *Bioinformatics* 35(10):1720–1728.
18. 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.
19. 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.
20. Ardizzone L, et al. (2018) Analyzing inverse problems with invertible neural networks. *arXiv preprint arXiv:1808.04730*.
21. Kingma DP, Dhariwal P (2018) Glow: Generative flow with invertible 1x1 convolutions in *Advances in Neural Information Processing Systems*. pp. 10215–10224.
22. 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*.
23. Dinh L, Sohl-Dickstein J, Bengio S (2016) Density estimation using real nvp. *arXiv preprint arXiv:1605.08803*.
24. Bloem-Reddy B, Teh YW (2019) Probabilistic symmetry and invariant neural networks. *arXiv preprint arXiv:1901.06082*.
25. Goodfellow I, Bengio Y, Courville A (2016) *Deep learning*. (MIT press).
26. 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.
27. Gelman A, et al. (2013) *Bayesian data analysis*. (Chapman and Hall/CRC).
28. 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.
29. Kingma DP, Welling M (2014) Auto-encoding variational bayes. *stat* 1050:1.