

---

# Neural Empirical Bayes: Source Distribution Estimation and its Applications to Simulation-Based Inference

---


Maxime Vandegar  
SLAC National  
Accelerator Laboratory

Michael Kagan  
SLAC National  
Accelerator Laboratory

Antoine Wehenkel  
ULiège

Gilles Louppe  
ULiège

## Abstract

We revisit empirical Bayes in the absence of a tractable likelihood function, as is typical in scientific domains relying on computer simulations. We investigate how the empirical Bayesian can make use of neural density estimators first to use all noise-corrupted observations to estimate a prior or source distribution over uncorrupted samples, and then to perform single-observation posterior inference using the fitted source distribution. We propose an approach based on the direct maximization of the log-marginal likelihood of the observations, examining both biased and de-biased estimators, and comparing to variational approaches. We find that, up to symmetries, a neural empirical Bayes approach recovers ground truth source distributions. With the learned source distribution in hand, we show the applicability to likelihood-free inference and examine the quality of the resulting posterior estimates. Finally, we demonstrate the applicability of Neural Empirical Bayes on an inverse problem from collider physics. 

## 1 Introduction

The estimation of a *source* distribution over latent random variables  $\mathbf{x}$  which give rise to a set of observations  $\mathbf{y}$ , after undergoing a potentially non-linear corruption process (i.e., a pushforward), is an inverse problem frequently of interest to the scientific and engineering communities. The source distribution,  $p(\mathbf{x})$ , may represent the distribution of plausible measure-

ments, or intermediate random variables in a hierarchical model, prior to corruption by a measurement or detection apparatus. The source distribution is of scientific interest as it allows comparison with theoretical predictions and for posterior inference for subsequent observations. Notably, in many scientific domains, the relationship between the source and observed distributions is encoded in a simulator that provides an approximation of the corruption process and generates samples from the likelihood  $p(\mathbf{y}|\mathbf{x})$ . However, as is typical with computer simulations, the likelihood function is implicit and rarely known in a tractable closed form.

Formally, we state the problem of likelihood-free source estimation as follows. Given a first dataset  $\mathbf{Y} = \{\mathbf{y}_i\}_{i=1}^N$  of  $N$  noise-corrupted observations  $\mathbf{y}_i$  and a second dataset  $\{\mathbf{x}_j, \mathbf{y}_j\}_{j=1}^M$  of matching pairs of source data and observations, with  $\mathbf{x}_j \sim \pi(\mathbf{x})$  drawn from an arbitrary proposal distribution  $\pi(\mathbf{x})$  and  $\mathbf{y}_j \sim p(\mathbf{y}|\mathbf{x}_j)$ , our aim is to learn the source distribution  $p(\mathbf{x})$ , not necessarily equal to  $\pi(\mathbf{x})$ , that has generated the observations  $\mathbf{Y}$ . For the class of problems we consider, we may assume that the dataset of  $(\mathbf{x}_j, \mathbf{y}_j)$  pairs is generated beforehand using a simulator of the stochastic corruption process.

The source distribution estimation problem is closely related to likelihood-free inference (LFI, [Cranmer et al., 2020](#)), though there are notable differences in problem statements. First, in Bayesian LFI, the objective is the computation of a posterior given a known prior and an implicit likelihood function. In our problem statement, the primary objective is rather to identify an unknown prior or source distribution that, once identified, then enables likelihood-free inference. Second, we only assume access to a pre-generated dataset of pairs of simulated source data and observations. In many settings, simulators are highly complex, with long run times to generate data. As such, sequential methods based on active calls to the simulator, as often found in the LFI literature, would be impractical.

In this work, we follow an empirical Bayes (EB, [Robbins, 1956](#); [Dempster et al., 1977](#)) approach to address

this challenge, using modern neural density estimators to approximate both the intractable likelihood and the unknown source distribution. Our method, which we call Neural Empirical Bayes (NEB), proceeds in two steps. First, using simulated pairs  $\{\mathbf{x}_j, \mathbf{y}_j\}_{j=1}^M$ , we use neural density estimation to learn an approximate likelihood. Second, by modeling the source distribution with a parameterized generative model, the log-marginal likelihood of the observations is approximated with Monte Carlo integration, and the parameters of the source distribution learned through gradient-based optimization. While our estimator of the log-marginal likelihood is biased, it is consistent and the use of deep generative models allows for fast and parallelizable Monte Carlo integration to mitigate its bias. Nonetheless, we also examine de-biased and variational estimators for comparison. Finally, once a source distribution and likelihood function are learned, we demonstrate that posterior inference for new observations may be performed with suitable sampling-based methods.

We first review EB and describe our NEB approach in Section 2, followed by an examination of log-marginal likelihood estimators in Section 3. Related work is discussed in Section 4. In Section 5, we present benchmark problems that explore the efficacy of NEB and provide comparison baselines, as well as a demonstration on a real-world application to collider physics. Further discussion and a summary are in Section 6. In addition, we provide a summary of the notations in Appendix A.

## 2 Empirical Bayes

Methods for EB (Robbins, 1956; Dempster et al., 1977) are usually divided into two estimation strategies (Efron, 2014): either modeling on the  $\mathbf{x}$ -space, called  $g$ -modeling; or on the  $\mathbf{y}$ -space, called  $f$ -modeling.

Here, we revisit  $g$ -modeling to learn a source distribution that regenerates the observations  $\mathbf{Y}$ . Specifically, we parameterize the source distribution as  $q_\theta(\mathbf{x})$  which, when passed through the likelihood  $p(\mathbf{y}|\mathbf{x})$ , results in a distribution  $q_\theta(\mathbf{y})$  over noisy observations. The log-marginal likelihood of the observations  $\mathbf{Y}$  is expressed as

$$\begin{aligned} \log q_\theta(\mathbf{Y}) &= \sum_{i=1}^N \log q_\theta(\mathbf{y}_i) \\ &= \sum_{i=1}^N \log \int p(\mathbf{y}_i|\mathbf{x}) q_\theta(\mathbf{x}) d\mathbf{x}, \end{aligned} \quad (1)$$

and its direct maximization with respect to the parameters  $\theta$  leads to a solution for the source distribution.

The maximization of the log-marginal likelihood is equivalent to the minimization of the Kullback–Leibler divergence  $\text{KL}(p(\mathbf{y})||q_\theta(\mathbf{y})) = \mathbb{E}_{p(\mathbf{y})} [-\log q_\theta(\mathbf{y})] + \mathcal{C} \approx -\frac{1}{N} \sum_{i=1}^N \log q_\theta(\mathbf{y}_i)$ . Therefore, as  $\mathbf{Y}$  increases, an optimal solution will correspond to a source distribution that exactly reproduces the observed distribution when passed through the corruption process. We note however that the maximization of Eq. 1 is an ill-posed problem: distinct source distributions may result in the same distribution over observations when folded through the corruption process. As a result, the learned source distribution may differ from the ground truth, for instance missing modes, but still reproduce the observed distribution. We discuss approaches to mitigate these undesired behaviors effects when *a priori* known properties of the source distribution are available.

In the likelihood-free setting, the likelihood function  $p(\mathbf{y}|\mathbf{x})$  is only implicitly defined by the simulator which prevents the direct estimation of Eq. 1. However, a dataset  $\{\mathbf{x}_j, \mathbf{y}_j\}_{j=1}^M$  can be generated beforehand by drawing uncorrupted samples  $\mathbf{x}_j$  from a proposal distribution  $\pi(\mathbf{x})$  and running the simulator to generate corresponding noise-corrupted observations  $\mathbf{y}_j$ . Similarly to Diggle and Gratton (1984) and D’Agostini (1995) who built likelihood function estimators with kernels or histograms, we use the generated dataset to train a surrogate  $q_\phi(\mathbf{y}|\mathbf{x})$  of the likelihood function, but we make use of modern neural density estimators such as normalizing flows (Tabak et al., 2010; Rezende and Mohamed, 2015). After the upfront simulation cost of generating the training data, no additional call to the simulator is needed.

We optimize the parameters  $\phi$  by maximizing the total log-likelihood  $\sum_{m=1}^M \log q_\phi(\mathbf{y}_m|\mathbf{x}_m)$  with mini-batch stochastic gradient ascent. Again, for large  $M$ , this is equivalent to minimizing  $\mathbb{E}_{\pi(\mathbf{x})} \text{KL}(p(\mathbf{y}|\mathbf{x})||q_\phi(\mathbf{y}|\mathbf{x}))$  and given enough capacity the surrogate likelihood is guaranteed to be a good approximation of  $p(\mathbf{y}|\mathbf{x})$  in the support of the proposal distribution  $\pi(\mathbf{x})$ . As a consequence, the support of  $\pi(\mathbf{x})$  should be chosen to cover the full range of plausible source data values. For example in the simulation-based inference setting,  $\pi(\mathbf{x})$  may be the distribution obtained from a simulator.

## 3 Log-marginal likelihood estimation

In Section 3.1 (resp. 3.2) we build a biased (resp. unbiased) estimator of the log-marginal likelihood  $\log q_\theta(\mathbf{y})$  with a generative model  $\mathbf{G}_\theta(\cdot) : \mathcal{E} \rightarrow \mathcal{X}$  that defines a differentiable mapping from a base distribution  $p(\epsilon)$  to  $q_\theta(\mathbf{x})$ . Then, in Section 3.3, we show how to use variational estimators of  $\log q_\theta(\mathbf{y})$  for NEB.

### 3.1 Biased estimator

Given a likelihood function  $p(\mathbf{y}|\mathbf{x})$  or its surrogate  $q_\phi(\mathbf{y}|\mathbf{x})$  we define an estimator  $\mathcal{L}_K(\theta)$  of the log-marginal likelihood. This estimator can be plugged in Eq. 1 to optimize the source distribution parameters  $\theta$  by stochastic minibatch gradient ascent. Based on Monte Carlo integration, the estimator is defined as:

$$\begin{aligned} \log q_\theta(\mathbf{y}) &= \log \mathbb{E}_{q_\theta(\mathbf{x})} [p(\mathbf{y}|\mathbf{x})] \\ &= \log \mathbb{E}_{p(\epsilon)} [p(\mathbf{y}|\mathbf{G}_\theta(\epsilon))] \\ &\approx \log \frac{1}{K} \sum_{k=1}^K p(\mathbf{y}|\mathbf{G}_\theta(\epsilon_k)) \\ &= \log \text{SumExp} [\log p(\mathbf{y}|\mathbf{G}_\theta(\epsilon_k))] - C \\ &=: \mathcal{L}_K(\theta), \end{aligned} \quad (2)$$

where  $\epsilon_k \sim p(\epsilon)$ ,  $C$  is a constant independent of  $\theta$ , and the log-sum-exp trick is used for numerical stability. While a large number  $K$  of samples may be needed for good Monte Carlo approximation, this difficulty is alleviated by the ease of generating large samples of source data with the neural sampler  $\mathbf{G}_\theta$ .

We study and prove properties of the estimator  $\mathcal{L}_K(\theta)$  in Appendix B. Using the Jensen's inequality, we first show that  $\mathcal{L}_K(\theta)$  is biased. We demonstrate however that both its bias and its variance decrease at a rate of  $\mathcal{O}(\frac{1}{K})$ . Then, similarly to Burda et al. (2016), we show that  $\mathcal{L}_K(\theta)$  is monotonically non-decreasing in expectation with respect to  $K$ , i.e.

$$\mathbb{E} [\mathcal{L}_{K+1}(\theta)] \geq \mathbb{E} [\mathcal{L}_K(\theta)]. \quad (3)$$

As  $K \rightarrow \infty$ , we finally show that the estimator is however consistent:

$$\lim_{K \rightarrow \infty} \mathcal{L}_K(\theta) = \log q_\theta(\mathbf{y}). \quad (4)$$

### 3.2 Unbiased estimator

Using the Russian roulette estimator (Kahn, 1955), we de-bias the log-marginal likelihood estimator  $\mathcal{L}_K(\theta)$  as

$$\hat{\mathcal{L}}_K(\theta) := \mathcal{L}_K(\theta) + \eta(\theta), \quad (5)$$

where  $\eta(\theta)$  is a random variable – whose expectation corrects for the bias – defined as

$$\eta(\theta) = \sum_{j=0}^J \frac{\mathcal{L}_{K+j+1}(\theta) - \mathcal{L}_{K+j}(\theta)}{P(\mathcal{J} \geq j)}, \quad (6)$$

with  $J \sim P(J)$ . Similarly to Luo et al. (2020) in their study of Importance Weighted Auto-Encoders (IWAEs, Burda et al., 2016), we prove in Appendix B.5

that  $\hat{\mathcal{L}}_K$  is an unbiased estimator as long as  $P(J)$  is a discrete distribution such that  $P(\mathcal{J} \geq j) > 0, \forall j > 0$ . Ideally, the distribution  $P(J)$  should be chosen such that it adds only a small computational overhead, while providing a finite-variance estimator. In our experiments, we reduce the computational overhead by re-using the same Monte Carlo terms used for  $\mathcal{L}_{K+j}$  to compute  $\mathcal{L}_{K+j+1}$ .

### 3.3 Variational empirical Bayes

For EB in high-dimension, Wang et al. (2019) proposed to build upon Kingma and Welling (2014) and to introduce a variational posterior distribution  $q_\psi(\mathbf{x}|\mathbf{y})$  whose parameters  $\psi$  are jointly optimized with the parameters  $\theta$  of the source distribution by maximizing the evidence lower bound (ELBO):

$$\begin{aligned} \log q_\theta(\mathbf{y}) &\geq \log q_\theta(\mathbf{y}) - \text{KL}(q_\psi(\mathbf{x}|\mathbf{y})||p(\mathbf{x}|\mathbf{y})) \\ &= \mathbb{E}_{q_\psi(\mathbf{x}|\mathbf{y})} [\log p(\mathbf{y}|\mathbf{x})] \\ &\quad - \text{KL}(q_\psi(\mathbf{x}|\mathbf{y})||q_\theta(\mathbf{x})) \\ &=: \mathcal{L}^{\text{ELBO}}. \end{aligned} \quad (7)$$

When  $\mathcal{L}^{\text{ELBO}}$  is optimized with stochastic gradient descent, an unbiased estimator can be obtained with Monte Carlo integration – usually only one Monte Carlo sample is used which yields a tractable objective. While being tractable, the ELBO is a lower bound (and biased estimator) of the log-marginal likelihood. A common approach (Rezende and Mohamed, 2015) to tighten the bound is to model  $q_\psi(\mathbf{x}|\mathbf{y})$  from a large distribution family so that it can closely match the posterior distribution, i.e. efficiently minimize  $\text{KL}(q_\psi(\mathbf{x}|\mathbf{y})||p(\mathbf{x}|\mathbf{y}))$ . Close to our work, IWAEs trade off computational complexity to obtain a tighter log-likelihood lower bound derived from importance sampling. Specifically, IWAEs are trained to maximize

$$\mathcal{L}_K^{\text{IW}}(\theta, \psi) = \log \frac{1}{K} \sum_{k=1}^K p(\mathbf{y}|\mathbf{x}_k) w(\mathbf{x}_k), \quad (8)$$

where  $w(\mathbf{x}_k) = \frac{q_\theta(\mathbf{x}_k)}{q_\psi(\mathbf{x}_k|\mathbf{y})}$  and  $\mathbf{x}_k \sim q_\psi(\mathbf{x}|\mathbf{y})$ . IWAEs are a generalization of the ELBO based on importance weighting (setting  $K = 1$  retrieves the ELBO objective). Nowozin (2018) showed that the bias and variance of this estimator vanish for  $K \rightarrow \infty$  at the same rate  $\mathcal{O}(\frac{1}{K})$  as  $\mathcal{L}_K$ .

By design,  $\mathcal{L}^{\text{ELBO}}$  and  $\mathcal{L}^{\text{IW}}$  require the evaluation of the density of new data points under the source model, whereas  $\mathcal{L}_K$  and  $\hat{\mathcal{L}}_K$  only require the efficient sampling from the source model. Which method to use should therefore depend on the downstream usage of the source distribution. While the evaluation of densities required by  $\mathcal{L}^{\text{ELBO}}$  and  $\mathcal{L}^{\text{IW}}$  limits the range of

models that can be used and makes the introduction of inductive bias more difficult,  $\mathcal{L}_K$  and  $\hat{\mathcal{L}}_K$  can be used with any generative model. In the rest of the manuscript we refer indistinguishably to the source distribution as  $q_\theta(\mathbf{x})$  although the generative models used with  $\mathcal{L}_K$  and  $\hat{\mathcal{L}}_K$  may not allow its evaluation. In that case, we mean the pushforward distribution implied by  $\mathbf{G}_\theta(\mathbf{x})$ .

## 4 Related work

**Empirical Bayes** In the most common forms of  $g$ -modeling, the likelihood function and the prior distribution are chosen such that the marginal likelihood can be computed and maximized iteratively or analytically. More recent approaches model the prior distribution analytically but assume both the  $\mathbf{x}$ -space and  $\mathbf{y}$ -space are finite and discrete (Narasimhan and Efron, 2016; Efron, 2016). Then, given a known likelihood function encoded in tensor form, the distribution parameters are optimized by maximum marginal likelihood estimation. Similarly to this latter approach, we do not require a likelihood function in closed-form, but we build a continuous surrogate that allows its direct evaluation rather than discretizing it.

While Wang et al. (2019) only theoretically proposed using Eq. 7 in EB, we show experimentally in the next section the applicability of this method. Concurrent work (Dockhorn et al., 2020) also used this approach to solve a density deconvolution task on Gaussian noise processes. Our work differs as we show the applicability of these methods on much more complicated black-box simulators, including a real inverse problem from collider physics. Black-box simulators imply that a neural network surrogate replaces the likelihood function, and thus, learning  $\theta$  and  $\psi$  requires to backpropagate through the surrogate.

Finally, in the context of likelihood-free inference, Louppe et al. (2019) used adversarial training for learning a prior distribution such that, when corrupted by a non-differentiable black-box model, reproduces the empirical distribution of the observations. This can be seen as  $g$ -modeling EB where a prior distribution is optimized based on observations.

**Unfolding** Approximating a source distribution  $p(\mathbf{x})$  given corrupted observations is often referred to as unfolding in the particle physics literature (for reviews see Cowan, 2002; Blobel, 2011; Adye, 2011). A common approach (Richardson, 1972; Lucy, 1974; D’Agostini, 1995) is to discretize the problem and replace the integral in Eq. 1 with a sum, resulting in a discrete linear inverse problem. The surrogate model  $q_\phi(\mathbf{y}|\mathbf{x})$  of the likelihood function is encoded in tensor

form while  $q_\theta(\mathbf{x})$  is modeled with a histogram. These approaches are typically limited to low dimensions. In order to scale to higher dimensions, preliminary work by Cranmer (2018) explored the idea of modelling  $q_\phi(\mathbf{y}|\mathbf{x})$  and  $q_\theta(\mathbf{x})$  with normalizing flows to approximate the integral in Eq. 1 with Monte Carlo integration. Aiming to the same objective, Andreassen et al. (2019b) replaced the sum in discrete space with a full-space integral using the likelihood ratio which is used for re-weighting. Bellagente et al. (2020) used invertible neural networks for learning a posterior that can be used for unfolding while our EB approach focuses on learning a source distribution at inference time.

**Likelihood-free inference** The use of a surrogate model of the likelihood function that enables inference as if the likelihood was known is not new. Since Diggle and Gratton (1984), kernels and histograms have been vastly used for 1D density estimation. More recently, several Bayesian likelihood-free inference algorithms (Papamakarios et al., 2019; Papamakarios and Murray, 2016; Lueckmann et al., 2017; Greenberg et al., 2019; Hermans et al., 2020; Durkan et al., 2020) have been developed to carry out inference when the likelihood function is implicit and intractable. These methods all operate by learning parts of the Bayes’ rule, such as the likelihood function, the likelihood-to-evidence ratio, or the posterior itself, and all require the explicit specification of a prior distribution. By contrast, the primary objective of our work is to learn a prior distribution from a set of noise-corrupted observations which, once it is identified, then enables any of the aforementioned Bayesian LFI algorithms for posterior inference. We refer the reader to Cranmer et al. (2020) for a broader review of likelihood-free inference.

## 5 Experiments

We present three studies of NEB. In Section 5.1, we analyze the intrinsic quality of the recovered source distribution for the estimators discussed in Section 3. In Section 5.2, we explore posterior inference with the learned source distribution. Finally, in Section 5.3 we show the applicability of NEB on an inverse problem from collider physics. All experiments are repeated 5 times with  $q_\phi(\mathbf{y}|\mathbf{x})$  and  $q_\theta(\mathbf{x})$  relearned in each experiment. Means and standard deviations are reported.

### 5.1 Source estimation

We evaluate NEB on three benchmark problems: (a) a toy model with a simple likelihood but complex posterior (SLCP) introduced by Papamakarios et al. (2019), (b) the two-moons model of Greenberg et al. (2019), and (c) an inverse kinematics problem (IK) proposed



Simulator	K	$\mathcal{L}_K$	$\hat{\mathcal{L}}_K$
SLCP	10	$0.82 \pm 0.01$	$0.65 \pm 0.04$
	128	$0.57 \pm 0.01$	$0.59 \pm 0.02$
	256	$0.55 \pm 0.02$	$0.54 \pm 0.00$
	1024	$0.53 \pm 0.01$	$0.52 \pm 0.01$
Two-moons	10	$0.69 \pm 0.02$	$0.56 \pm 0.02$
	128	$0.53 \pm 0.01$	$0.57 \pm 0.06$
	256	$0.52 \pm 0.01$	$0.52 \pm 0.02$
	1024	$0.52 \pm 0.01$	$0.53 \pm 0.01$
IK	10	$0.80 \pm 0.13$	$0.67 \pm 0.08$
	128	$0.65 \pm 0.04$	$0.67 \pm 0.12$
	256	$0.66 \pm 0.02$	$0.71 \pm 0.09$
	1024	$0.66 \pm 0.03$	$0.62 \pm 0.03$

Table 1: ROC AUC between the observed distribution  $p(\mathbf{y})$  and the regenerated distribution  $\int p(\mathbf{y}|\mathbf{x})q_\theta(\mathbf{x})d\mathbf{x}$ . The closer to 0.5, the better the estimation in  $\mathbf{y}$ -space. *NEB successfully identifies source distributions that result in distributions over noise-corrupted observations that are almost indistinguishable from the ground truth. When  $K$  is low, de-biasing leads to substantial improvements.*

by Ardizzone et al. (2018). See Appendix C for complementary experimental details. We use datasets of  $M = 15000$  samples to train surrogate models  $q_\phi(\mathbf{y}|\mathbf{x})$  for each simulator. All density models are parameterized with normalizing flows made of four coupling layers (Dinh et al., 2015, 2017). Further architecture and optimization details can be found in Appendix D. The source distributions  $q_\theta(\mathbf{x})$  are optimized on  $N = 10000$  observations  $\mathbf{y}$  and we show further results with only one or two observations in Appendix E. The ground truth source distributions  $p(\mathbf{x})$  are  $\mathcal{U}(-3, 3)^5$  for SLCP,  $\mathcal{U}(-1, 1)^2$  for two-moons and  $\mathcal{N}(\mathbf{0}, \text{Diag}(\frac{1}{4}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2}))$  for IK.

**Biased vs. unbiased estimator** We first compare the biased and unbiased estimators  $\mathcal{L}_K$  and  $\hat{\mathcal{L}}_K$ . Table 1 reports the ROC AUC scores of a classifier trained to distinguish between noise-corrupted observations from the ground truth  $p(\mathbf{y})$  and noise-corrupted observations from the marginal  $\int p(\mathbf{y}|\mathbf{x})q_\theta(\mathbf{x})d\mathbf{x}$  obtained by passing source data from  $q_\theta(\mathbf{x})$  into the exact simulator. For both estimators, the table shows that we successfully identify a source distribution  $q_\theta(\mathbf{x})$  resulting in a distribution over noise-corrupted observations which is almost indistinguishable from the ground truth  $p(\mathbf{y})$ . When  $K$  is low, de-biasing the estimator leads to significant improvements. When  $K$  increases, the bias of  $\mathcal{L}_K$  drops quickly, and de-biasing, which introduces variance, does not significantly improve the results. Therefore, we recommend using the de-biased estimator when  $K$

is constrained to be low, e.g., when the GPU memory is limited. In the following, we set  $K = 1024$  and only consider the biased estimator  $\mathcal{L}_K(\theta)$ .

**Monte Carlo vs. variational methods** We evaluate the quality of  $\mathcal{L}_K$ ,  $\mathcal{L}^{\text{ELBO}}$  and  $\mathcal{L}_K^{\text{IW}}$ . For a fair comparison, we use an Unconstrained Monotonic Neural Network autoregressive flow (UMNN-MAF, Wehenkel and Louppe, 2019) to parameterize the prior for all losses. The recognition network  $q_\psi(\mathbf{x}|\mathbf{y})$  for the variational approaches is modeled with the same architecture as the prior, but is conditioned on  $\mathbf{y}$ . We use  $K = 128$  for  $\mathcal{L}_K^{\text{IW}}$  due to GPU memory constraints. We show in Appendix F that simpler implicit generative models can be used with the Monte Carlo estimators  $\mathcal{L}_K$  and  $\hat{\mathcal{L}}_K$ , effectively reducing inference time and allowing to use higher values of  $K$ .

Table 2 shows the ROC AUC of a classifier trained to discriminate between samples from the ground truth source distribution  $p(\mathbf{x})$  and samples from the source distribution  $q_\theta(\mathbf{x})$  identified by each of the different methods. A ROC AUC score between 0.5 and 0.7 is often considered poor discriminative performance, therefore indicating good source estimation. The estimator  $\mathcal{L}_{1024}$  leads to the most accurate source distributions on these three tasks. In particular, the source distribution found for the two-moons problem is almost perfect. At the same time, the results for SLCP and IK are marginally acceptable, and largely better than for the variational methods ( $\mathcal{L}_K^{\text{IW}}$  and  $\mathcal{L}^{\text{ELBO}}$ ). Figures 1 and 2 illustrate for  $\mathcal{L}_{1024}$  how the exact and learned sources distributions are visually similar.

Table 2 also reports the discrepancy between the corrupted data from the identified source distributions and the ground truth distribution of noise-corrupted observations. While  $\mathcal{L}^{\text{ELBO}}$  does not give good results for SLCP, tightening the evidence lower-bound with  $\mathcal{L}_{128}^{\text{IW}}$  yields good results on all problems. While  $\mathcal{L}_{1024}$  has similar performance to  $\mathcal{L}_{128}^{\text{IW}}$  on SLCP and two-moons, it is performing worse for IK, due to the difficulty of approximating  $\int p(\mathbf{y}|\mathbf{x})p(\mathbf{x})d\mathbf{x}$  from Monte Carlo integration since the likelihood function for this problem is almost a Dirac function (see Appendix C for more details).

After observing the different estimators’ reconstruction quality, the superiority of  $\mathcal{L}_{1024}$  on source estimation over variational methods may be surprising at first glance. However, the variational methods require learning both a source distribution and a recognition network that are consistent with the likelihood function and the observations. This means that a wrong recognition network may prevent learning the correct source distribution as they must be consistent with each other. In the three experiments analyzed here,

Simulator	x-space			y-space			x-space (symmetric prior)		
	$\mathcal{L}^{\text{ELBO}}$	$\mathcal{L}^{\text{IW}}_{128}$	$\mathcal{L}_{1024}$	$\mathcal{L}^{\text{ELBO}}$	$\mathcal{L}^{\text{IW}}_{128}$	$\mathcal{L}_{1024}$	$\mathcal{L}^{\text{ELBO}}$	$\mathcal{L}^{\text{IW}}_{128}$	$\mathcal{L}_{1024}$
SLCP	$1.00 \pm 0.00$	$0.82 \pm 0.09$	$0.75 \pm 0.03$	$0.92 \pm 0.04$	$0.50 \pm 0.00$	$0.53 \pm 0.01$	$0.99 \pm 0.01$	$0.59 \pm 0.05$	$0.81 \pm 0.02$
Two-Moons	$0.75 \pm 0.00$	$0.75 \pm 0.00$	$0.55 \pm 0.02$	$0.50 \pm 0.01$	$0.50 \pm 0.00$	$0.52 \pm 0.01$	$0.51 \pm 0.01$	$0.50 \pm 0.01$	$0.51 \pm 0.02$
IK	$1.00 \pm 0.00$	$0.95 \pm 0.05$	$0.74 \pm 0.03$	$0.51 \pm 0.01$	$0.50 \pm 0.01$	$0.62 \pm 0.03$	$0.97 \pm 0.01$	$0.72 \pm 0.02$	$0.66 \pm 0.04$

Table 2: Source estimation for the benchmark problems. ROC AUC between  $q_\theta(\mathbf{x})$  and  $p(\mathbf{x})$  (x-space), and between the observed distribution  $p(\mathbf{y})$  and the regenerated distribution  $\int p(\mathbf{y}|\mathbf{x})q_\theta(\mathbf{x})d\mathbf{x}$  (y-space).

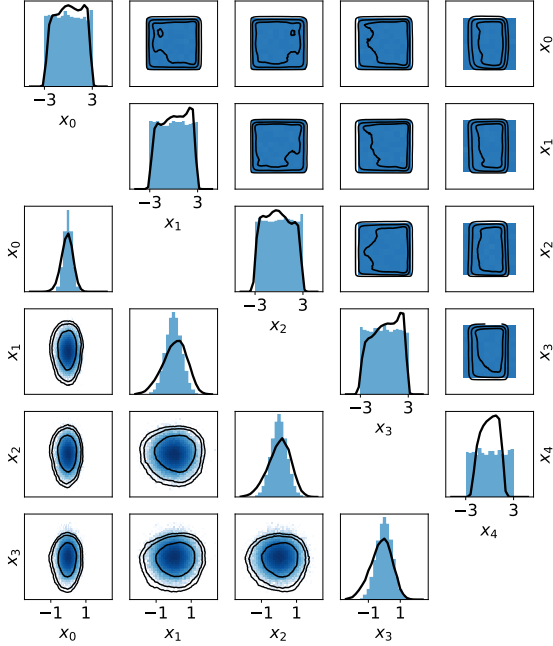


Figure 1: Source estimation results for  $\mathcal{L}_{1024}$  on SLCP (top) and IK (bottom). The source distribution  $p(\mathbf{x})$  are shown in blue against the estimated source distribution  $q_\theta(\mathbf{x})$  in black (the 68-95-99.7% contours are shown). The identified source distributions are similar to the unseen source distributions.

the prior distribution is a simple unimodal continuous distribution, whereas the posteriors are discontinuous and multimodal. In these cases, learning only the source distribution is simpler than learning both a source and posterior distributions.

**Symmetric source distribution** As mentioned before, multiple optimal solutions may co-exist when the inverse problem is ill-posed. On close inspection, figures 1 and 2 show that NEB successfully recovers the domain of the source data but fails to exactly reproduce the ground truth source distribution. Indeed, for all problems considered here, the passage of  $\mathbf{x}$  through the corruption process results in a loss of information in  $\mathbf{y}$ , which may lead to multiple solutions. We observe this in Figure 2, where we plot the quantities  $|\mathbf{x}_1 + \mathbf{x}_2|$

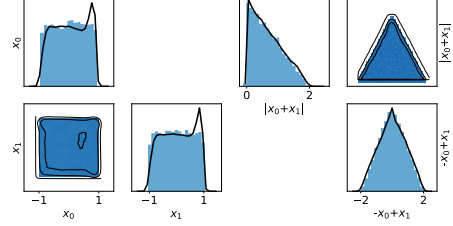


Figure 2: Source estimation results for  $\mathcal{L}_{1024}$  on the two-moons problem. The source distribution  $p(\mathbf{x})$  is shown in blue against the estimated source distribution  $q_\theta(\mathbf{x})$  in black (the 68-95-99.7% contours are shown). As shown on the right, up to the symmetries of the problem, the identified source distribution matches the unseen source distribution.

and  $-\mathbf{x}_1 + \mathbf{x}_2$  that are sufficient statistics of  $\mathbf{x}$  for estimating  $\mathbf{y}$ . We see that the distribution over these intermediate variables is nearly equal for the ground truth distribution and the identified source distribution. This indicates that, up to symmetries, NEB recovers the ground truth source distribution.

A reasonable way to encourage learning a good source distribution is to enforce *a priori* known properties such as its domain, symmetries, or smoothness. This type of useful inductive biases can be embedded in the neural network to constrain the solution space. As such, we modify the UMMN-MAF networks  $q_\theta(\mathbf{x})$  so that the generated distributions are one-to-one symmetric, i.e.  $q_\theta([x_1, \dots, x_d]) = q_\theta([\pm x_1, \dots, \pm x_d])$  (these modifications are detailed in Appendix G). Table 2 shows that the symmetric distributions are more similar to the unseen distributions  $p(\mathbf{x})$  in all but one case. For example, all methods learn to approximately identify the exact source distribution on the two-moons despite the simulator’s destructive process. Results for  $\mathcal{L}_{1024}$  on SLCP are worse because the regularization pushes the learned distribution to a solution that still reproduces the observed distribution with high accuracy (the ROC AUC between the observed distribution and the regenerated one drops to  $0.51 \pm 0.01$ ), but that moves away from the unseen source distribution. Further inductive bias should therefore be introduced; for example, the learned distribution can be bounded

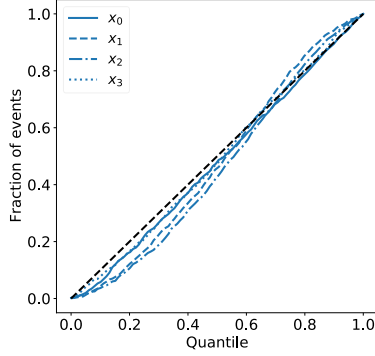


Figure 3: Posterior inference for IK. The plot shows the per-parameter calibration curves obtained with rejection sampling using the learned likelihood  $q_\psi(\mathbf{y}|\mathbf{x})$  and the identified source distribution  $q_\theta(\mathbf{x})$  with  $\mathcal{L}_{1024}$ . The curves indicate a reasonably well calibrated posterior distribution.

using specific activation functions in the last layer of  $\mathbf{G}_\theta(\cdot)$ . We note here that the non-variational methods are generally better suited for inductive bias as their architecture designs are less constrained.

## 5.2 Likelihood-free posterior inference

In the context of Bayesian posterior inference, the source distribution we retrieve with NEB can be used as a prior distribution. Therefore, the learned prior  $q_\theta(\mathbf{x})$  together with the surrogate likelihood  $q_\phi(\mathbf{y}|\mathbf{x})$  unlock the subsequent likelihood-free estimation of the posterior  $p(\mathbf{x}|\mathbf{y})$  – for which the fidelity will depend on both the correctness of the source distribution and the likelihood. There is an ongoing debate in the EB literature regarding the use of the data twice for posterior inference in this approach (Gelman, 2008; Darnieder, 2011; Gelman et al., 2017). When few prior knowledge are available, EB allows to learn insights from the data, and we believe it is valuable to incorporate those insights in the prior rather than, for instance, choosing a wide prior and especially in high dimensions where the prior choice is important (Gelman et al., 2017).

As we have used normalizing flows with  $\mathcal{L}_K$  and  $\hat{\mathcal{L}}_K$ , state-of-the-art Markov Chain Monte Carlo (MCMC) methods such as Hybrid Monte Carlo (HMC) can be used for sampling the posterior. Other generative models that do not allow density evaluation could be used in our empirical Bayes setup but would not permit the usage of MCMC. In this section, we focus on rejection sampling, rather than MCMC, as the source distribution model allows fast parallel sampling which makes the algorithm efficient even when the acceptance rate is low. We perform rejection sampling as follows: given  $u \sim \mathcal{U}(0, 1)$ , we accept samples  $\mathbf{x} \sim q_\theta(\mathbf{x})$

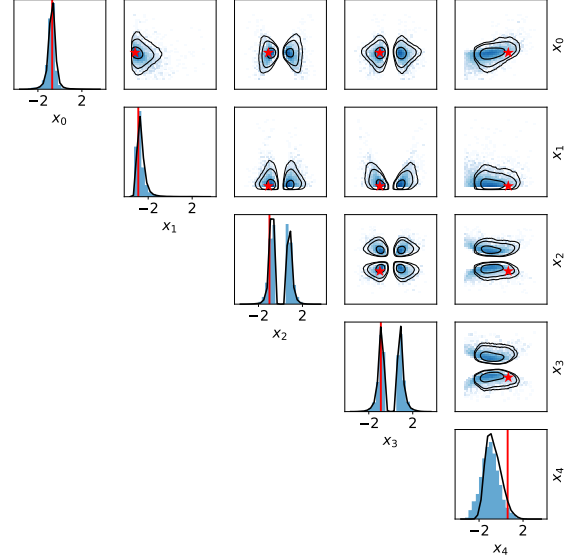


Figure 4: Posterior distribution obtained from MCMC with the exact source distribution and the exact likelihood function on SLCP in blue against the posterior distribution obtained with  $q_\phi(\mathbf{y}|\mathbf{x})$  and  $q_\theta(\mathbf{x})$  learned from  $\mathcal{L}_{1024}$  in black (the 68-95-99.7% contours are shown). Generating source sample  $\mathbf{x}$  are indicated in red. The approximated posterior distribution closely matches the ground truth.

such that  $u < \frac{q_\phi(\mathbf{y}|\mathbf{x})}{M}$  where  $M > q_\phi(\mathbf{y}|\mathbf{x}), \forall \mathbf{x}$  is determined empirically. On the other hand, variational approaches ( $\mathcal{L}^{\text{ELBO}}$  and  $\mathcal{L}_K^{\text{IW}}$ ) directly learn a posterior  $q_\psi(\mathbf{x}|\mathbf{y})$  as a function of the single observation  $\mathbf{y}$ , which enables immediate per-event posterior inference. For  $\mathcal{L}_K^{\text{IW}}$ , Cremer et al. (2017) suggest using importance sampling.

We assess the goodness of the posterior distributions with a calibration test. Inspired by Bellagente et al. (2020), for multiple observations  $\mathbf{y}_i$ , we approximate the 1D posterior distributions and we report the fraction of events as a function of the quantile to which the generating source data  $\mathbf{x}_i$  fall. Figure 3 reports calibration curves associated with  $\mathcal{L}_{1024}$  for the IK problem, indicating reasonably well calibrated posteriors. We should note however that this observation is a necessary but not sufficient condition for well-calibrated posterior distributions. More details and quantitative results are given in Appendix H.

Finally, we show in Figure 4 an example of posterior distribution obtained with rejection sampling using  $q_\theta(\mathbf{x})$  and  $q_\phi(\mathbf{y}|\mathbf{x})$  as learned with  $\mathcal{L}_{1024}$ , against the ground truth posterior obtained with Markov Chain Monte Carlo using the exact likelihood and source data distribution. We emphasize that NEB can recover nearly the exact posterior with no access to the

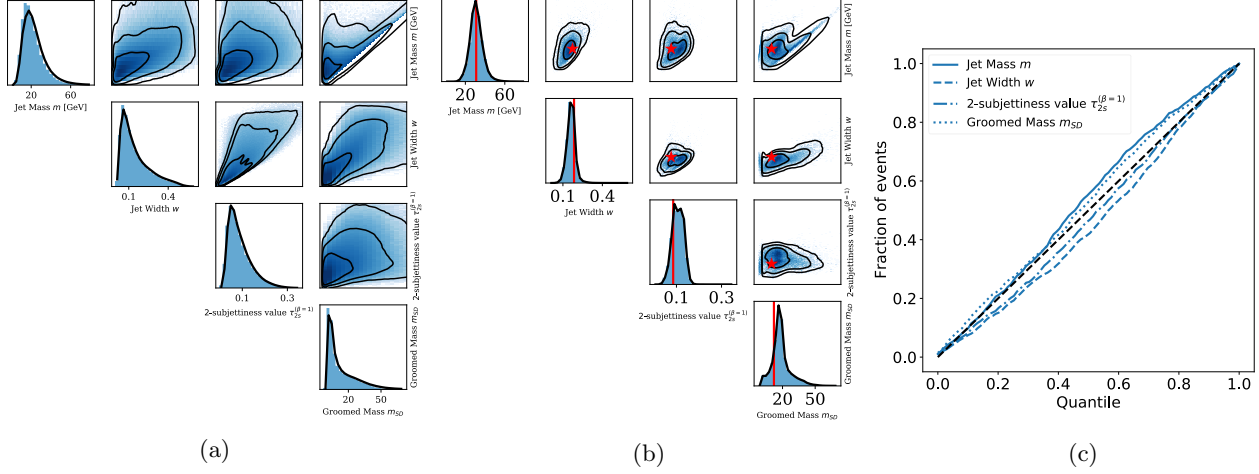


Figure 5: Neural Empirical Bayes for detector correction in collider physics. (a) The source distribution  $p(\mathbf{x})$  is shown in blue against the estimated source distribution  $q_\theta(\mathbf{x})$  in black. (b) Posterior distribution obtained with rejection sampling, with generating source sample  $\mathbf{x}$  indicated in red. (c) Calibration curves for each jet property obtained with rejection sampling on 10000 observations. In (a) and (b), contours represent the 68-95-99.7% levels.

likelihood function or to the prior distribution. To the best of our knowledge, this is the first work to show posterior inference is possible in this extreme setting.

### 5.3 Detector correction in collider physics

At colliders like the LHC, the distribution of particles produced from an interaction and incident on detectors can be predicted from theoretical models. Thus measurements of such distributions can be used to directly test theoretical predictions. However, while detectors measure the energy and momentum of particles, they also induce noise due to the stochastic nature of particle-material interactions and of the signal acquisition process. Thus a key challenge in comparing measurements to theoretical predictions is to correct noisy detector observations to obtain experimentally observed incident particle source distributions. This is frequently done by binning 1D or 2D distributions and solving a discrete linear inverse problem. Instead we apply NEB for estimating the multi-dimensional source distribution. We use the publicly available simulated dataset (Andreassen et al., 2019a) of paired source and corrupted measurements of properties of jets, or collimated streams of particles produced by high energy quarks and gluons. Simulation details are found in Appendix I.

Surrogate training was performed using one source simulator as a proposal distribution. The same surrogate architecture and hyperparameters as in the toy experiments were used (see Appendix D for details). We assess NEB in a dataset with the source distribution produced by a different simulator of the same

physical process. Both datasets for surrogate training and source distribution learning contain approximately 1.6 million events. This is an example setting where sequential inference methods cannot be used as only a fixed dataset is available and not the simulator.

	$\mathcal{L}^{\text{ELBO}}$	$\mathcal{L}_{128}^{\text{IW}}$	$\mathcal{L}_{1024}$
<b>x-space</b>	$0.99 \pm 0.02$	$0.63 \pm 0.06$	$0.57 \pm 0.05$
<b>y-space</b>	$0.87 \pm 0.08$	$0.51 \pm 0.01$	$0.50 \pm 0.01$

Table 3: Source estimation in collider physics. ROC AUC between  $q_\theta(\mathbf{x})$  and the unseen source distribution  $p(\mathbf{x})$  (x-space), and between the observed distribution  $p(\mathbf{y})$  and the regenerated distribution  $\int q_\phi(\mathbf{y}|\mathbf{x})q_\theta(\mathbf{x})d\mathbf{x}$ .

**Source estimation** We focus on the  $\mathcal{L}_{1024}$  estimator for source distribution learning although we also report results with the other estimators. Optimization is done with Adam using default parameters and an initial learning rate of  $10^{-4}$ . We train for 10 epochs with minibatches of size 256. The density estimator for the source distribution comprised 6 coupling layers, with 3-layer MLPs with 32 units per layer and ReLU activations used for the scaling and translation functions. Parameters were determined with a hyperparameter grid search using a held out validation set from the dataset on which the surrogate is trained. The learned source distribution is observed to closely match the true simulated source distribution, as seen in Figure 5a. Table 3 reports the ROC AUC between the learned and ground truth distributions, indicating only small discrepancies between them.



**Likelihood-free posterior inference** Figure 5b shows the learned posterior distribution against the generating source data. Plots are scaled to the prior-space. The model learns nicely a region of plausible values for the generating source data. To assess the quality of the posterior inference on more data, Figure 5c shows the fraction of events as a function of the quantile to which the generating source data belongs under the learned posterior distribution. Results indicate reasonably well calibrated posterior distributions.

## 6 Summary and discussion

In this work, we revisit  $g$ -modeling empirical Bayes with neural networks to estimate source distributions from non-linearly corrupted observations. We propose both a biased and de-biased estimator of the log-marginal likelihood, and examine variational methods for this challenge. We show that we can successfully recover source distributions from corrupted observations. We find that inductive bias is highly beneficial for solving ill-posed inverse problems and can be embedded in the structure of the neural networks used to model the source distribution. Although the explored approaches are general, we specifically study the likelihood-free setting, and we successfully perform posterior inference without direct access to either a likelihood function or a prior distribution.

**Future work** In this work we have mainly examined low-dimensional settings. We believe that further analysis of these methods for high-dimensional data such as images and time series could be of strong interest from both a theoretical and practical point of view. In particular, assessing the computational challenges of each method and the importance of inductive bias in this challenging setting are promising directions towards improvements in solving high-dimensional inverse problems.

## Acknowledgments

We thank Johann Brehmer and Kyle Cranmer for their helpful feedback on the manuscript. We also thank the anonymous reviewers for their thoughtful comments. Maxime Vandegar and Michael Kagan are supported by the US Department of Energy (DOE) under grant DE-AC02-76SF00515, and Michael Kagan is also supported by the SLAC Panofsky Fellowship. Antoine Wehenkel is a research fellow of the F.R.S.-FNRS (Belgium) and acknowledges its financial support. Gilles Louppe is recipient of the ULiège - NRB Chair on Big data and is thankful for the support of NRB.

## References

- ATLAS Pythia 8 tunes to 7 TeV datas. Technical Report ATL-PHYS-PUB-2014-021, CERN, Geneva, Nov 2014.
- Tim Adye. Unfolding algorithms and tests using RooUnfold. In *PHYSTAT 2011*, pages 313–318, Geneva, 2011. CERN. doi: 10.5170/CERN-2011-006.313.
- Anders Andreassen, Patrick Komiske, Eric Metodiev, Benjamin Nachman, and Jesse Thaler. Pythia/Herwig + Delphes Jet Datasets for Omnifold Unfolding, November 2019a.
- Anders Johan Andreassen, Patrick T. Komiske, Eric M. Metodiev, Benjamin Nachman, and Jesse Thaler. Omnifold: A method to simultaneously unfold all observables. 2019b.
- Jordanka A Angelova. On moments of sample mean and variance. *Int. J. Pure Appl. Math*, 79(1):67–85, 2012.
- Lynton Ardizzone, Jakob Kruse, Sebastian J. Wierker, Daniel Rahner, Eric W. Pellegrini, Ralf S. Klessen, Lena Maier-Hein, Carsten Rother, and Ullrich Köthe. Analyzing inverse problems with invertible neural networks. In *International Conference on Learning Representations*, 2018.
- Lynton Ardizzone, Carsten Lüth, Jakob Kruse, Carsten Rother, and Ullrich Köthe. Guided image generation with conditional invertible neural networks. *arXiv preprint arXiv:1907.02392*, 2019.
- ATLAS Collaboration. The ATLAS Experiment at the CERN Large Hadron Collider. *Journal of Instrumentation*, 3(08):S08003–S08003, aug 2008.
- Manuel Bahr, S. Gieseke, M. A. Gigg, David Grellscheid, K. Hamilton, O. Latunde-Dada, Simon Platzer, P. Richardson, Mike H. Seymour, A. Sherstnev, and Bryan Webber. Herwig++ 2.1 release note. 1999.
- Marco Bellagente, Anja Butter, Gregor Kasieczka, Tilman Plehn, Armand Rousselot, Ramon Winterhalder, Lynton Ardizzone, and Ullrich Köthe. Invertible networks or partons to detector and back again. *SciPost Physics*, 9(5):074, 2020.
- Volker Blobel. Unfolding Methods in Particle Physics. pages 240–251. 12 p, Jan 2011. doi: 10.5170/CERN-2011-006.240.
- Yuri Burda, Roger B Grosse, and Ruslan Salakhutdinov. Importance weighted autoencoders. In *International Conference on Learning Representations*, 2016.
- Manuel Bähr, Stefan Gieseke, Martyn A. Gigg, David Grellscheid, Keith Hamilton, Oluseyi Latunde-Dada, Simon Platzer, Peter Richardson, Michael H.

- Seymour, Alexander Sherstnev, and et al. Herwig++ physics and manual. 58(4):639–707, Nov 2008.
- Ricky TQ Chen, Jens Behrmann, David K Duvenaud, and Jörn-Henrik Jacobsen. Residual flows for invertible generative modeling. In *Advances in Neural Information Processing Systems*, pages 9916–9926, 2019.
- CMS Collaboration. The CMS experiment at the CERN LHC. *Journal of Instrumentation*, 3(08):S08004–S08004, aug 2008.
- Glen Cowan. A survey of unfolding methods for particle physics. *Conf. Proc. C*, 0203181:248–257, 2002.
- Kyle Cranmer. Neural unfolding. 2018. URL [https://github.com/cranmer/neural\\_unfolding](https://github.com/cranmer/neural_unfolding).
- Kyle Cranmer, Johann Brehmer, and Gilles Louppe. The frontier of simulation-based inference. *Proceedings of the National Academy of Sciences*, 2020.
- C. Cremer, Q. Morris, and D. Duvenaud. Reinterpreting Importance-Weighted Autoencoders. *Workshop at the International Conference on Learning Representations*, 2017.
- G. D’Agostini. A multidimensional unfolding method based on bayes’ theorem. *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment*, 362(2):487 – 498, 1995.
- William Francis Darnieder. *Bayesian methods for data-dependent priors*. PhD thesis, The Ohio State University, 2011.
- J. de Favereau, C. Delaere, P. Demin, A. Giammanco, V. Lemaître, A. Mertens, and M. Selvaggi. Delphes 3: a modular framework for fast simulation of a generic collider experiment. *Journal of High Energy Physics*, 2014(2), Feb 2014.
- A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22, 1977.
- Peter J. Diggle and Richard J. Gratton. Monte carlo methods of inference for implicit statistical models. 46(2):193–227, 1984.
- Laurent Dinh, David Krueger, and Yoshua Bengio. Nice: Non-linear independent components estimation. In *International Conference in Learning Representations workshop track*, 2015.
- Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp. In *International Conference in Learning Representations*, 2017.
- Tim Dockhorn, James A. Ritchie, Yaoliang Yu, and Iain Murray. Density deconvolution with normalizing flows. *Second workshop on Invertible Neural Networks, Normalizing Flows, and Explicit Likelihood Models (ICML 2020)*, 2020.
- Conor Durkan, Iain Murray, and George Papamakarios. On contrastive learning for likelihood-free inference. In *International Conference on Machine Learning*, 2020.
- Bradley Efron. Two modeling strategies for empirical bayes estimation. *Statistical science: a review journal of the Institute of Mathematical Statistics*, 29(2): 285, 2014.
- Bradley Efron. Empirical bayes deconvolution estimates. *Biometrika*, 103:1–20, 03 2016.
- Lyndon Evans and Philip Bryant. LHC machine. *Journal of Instrumentation*, 3(08):S08001–S08001, aug 2008.
- Andrew Gelman. Objections to bayesian statistics. *Bayesian Anal.*, 3(3):445–449, 09 2008. doi: 10.1214/08-BA318.
- Andrew Gelman, Daniel Simpson, and Michael Betancourt. The prior can often only be understood in the context of the likelihood. *Entropy*, 19(10):555, Oct 2017. ISSN 1099-4300. doi: 10.3390/e19100555.
- David Greenberg, Marcel Nonnenmacher, and Jakob Macke. Automatic posterior transformation for likelihood-free inference. In *International Conference on Machine Learning*, pages 2404–2414. PMLR, 2019.
- Joeri Hermans, Volodimir Begy, and Gilles Louppe. Likelihood-free mcmc with amortized approximate ratio estimators. In *International Conference on Machine Learning*, pages 4239–4248. PMLR, 2020.
- Herman Kahn. *Use of different Monte Carlo sampling techniques*. Rand Corporation, 1955.
- Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014.
- Gilles Louppe, Joeri Hermans, and Kyle Cranmer. Adversarial variational optimization of non-differentiable simulators. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1438–1447, 2019.
- Leon B Lucy. An iterative technique for the rectification of observed distributions. *The astronomical journal*, 79:745, 1974.
- Jan-Matthis Lueckmann, Pedro J. Goncalves, Giacomo Bassetto, Kaan Öcal, Marcel Nonnenmacher,

- and Jakob H. Macke. Flexible statistical inference for mechanistic models of neural dynamics. In *Advances in Neural Information Processing Systems*, 2017.
- Yucen Luo, Alex Beatson, Mohammad Norouzi, Jun Zhu, David Kristjanson Duvenaud, Ryan P. Adams, and Ricky T. Q. Chen. Sumo: Unbiased estimation of log marginal probability for latent variable models. *International Conference on Learning Representations.*, 2020.
- Balasubramanian Narasimhan and Bradley Efron. A g-modeling program for deconvolution and empirical bayes estimation. 2016.
- Sebastian Nowozin. Debiasing evidence approximations: On importance-weighted autoencoders and jackknife variational inference. In *International Conference on Learning Representations*, 2018.
- George Papamakarios and Iain Murray. Fast  $\epsilon$ -free inference of simulation models with bayesian conditional density estimation. In *Advances in Neural Information Processing Systems*. 2016.
- George Papamakarios, David Sterratt, and Iain Murray. Sequential neural likelihood: Fast likelihood-free inference with autoregressive flows. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 837–848. PMLR, 2019.
- Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International Conference on Machine Learning*, pages 1530–1538. PMLR, 2015.
- William Hadley Richardson. Bayesian-Based Iterative Method of Image Restoration. *Journal of the Optical Society of America (1917-1983)*, 62(1):55, January 1972.
- Herbert Robbins. An empirical bayes approach to statistics. In *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics*, pages 157–163, 1956.
- Torbjörn Sjöstrand, Stefan Ask, Jesper R. Christiansen, Richard Corke, Nishita Desai, Philip Ilten, Stephen Mrenna, Stefan Prestel, Christine O. Rasmussen, and Peter Z. Skands. An introduction to pythia 8.2. *Computer Physics Communications*, 191:159 – 177, 2015.
- Esteban G Tabak, Eric Vanden-Eijnden, et al. Density estimation by dual ascent of the log-likelihood. *Communications in Mathematical Sciences*, 8(1): 217–233, 2010.
- Yixin Wang, Andrew C Miller, David M Blei, et al. Comment: Variational autoencoders as empirical bayes. *Statistical Science*, 34(2):229–233, 2019.
- Antoine Wehenkel and Gilles Louppe. Unconstrained monotonic neural networks. In *Advances in Neural Information Processing Systems*, pages 1545–1555, 2019.