

---

# Likelihood-free inference with an improved cross-entropy estimator

---

Markus Stoye,<sup>1</sup> Johann Brehmer,<sup>2</sup> Gilles Louppe,<sup>3</sup> Juan Pavez,<sup>4</sup> and Kyle Cranmer<sup>2</sup>

<sup>1</sup> Department of Physics and Data Science Institute, Imperial College London

<sup>2</sup> Center for Cosmology and Particle Physics and Center for Data Science, New York University

<sup>3</sup> Department of Electrical Engineering and Computer Science, University of Liège,

<sup>4</sup> Federico Santa María Technical University

markus.stoye@cern.ch, johann.brehmer@nyu.edu, g.louppe@uliege.be,

juan.pavez@alumnos.usm.cl, kyle.cranmer@nyu.edu

## Abstract

We extend recent work (Brehmer, et. al., 2018) that use neural networks as surrogate models for likelihood-free inference. As in the previous work, we exploit the fact that the joint likelihood ratio and joint score, conditioned on both observed and latent variables, can often be extracted from an implicit generative model or simulator to augment the training data for these surrogate models. We show how this augmented training data can be used to provide a new cross-entropy estimator, which provides improved sample efficiency compared to previous loss functions exploiting this augmented training data.

## 1 Introduction

Many real-world phenomena are best described by computer simulations. Such simulators often implement a stochastic generative process, which is based on a mechanistic model and parametrized by  $\theta$ . In practice, these simulators are used to generate samples of observations  $x \sim p(x|\theta)$ , but the density is only defined implicitly through the simulation code. Often, the generative process involves latent variables and the density

$$p(x|\theta) = \int dz p(x, z|\theta) \quad (1)$$

is intractable because of the integral over a large (and possibly highly structured) latent space. Without a tractable likelihood, statistical inference on the parameters  $\theta$  given observed data  $x$  is challenging. This problem has prompted the development of *likelihood-free inference* methods such as Approximate Bayesian Computation [1–4] and neural density or neural density ratio estimation algorithms [5–23]. Nearly all of these established methods treat the simulator as a black box and only use its capability to generate samples for a specified values of  $\theta$ .

In Refs. [24–26] a new paradigm was introduced that exploits additional information that can be extracted from the simulation. In particular, within the simulation where the latent variables  $z$  are available, it is often possible to extract the *joint likelihood ratio*

$$r(x, z|\theta_0, \theta_1) = \frac{p(x, z|\theta_0)}{p(x, z|\theta_1)} \quad (2)$$

and the *joint score*

$$t(x, z|\theta_0) = \nabla_{\theta} \log p(x, z|\theta) \Big|_{\theta_0}, \quad (3)$$

which are conditioned on the latent variables  $z$  corresponding to a particular sample.

It was then shown that certain loss functionals  $L[g(x)]$ , which depend on the joint likelihood ratio and the joint score, are minimized by the likelihood ratio

$$g^*(x) \equiv \arg \min_{g(x)} L[g(x)] = r(x|\theta_0, \theta_1) \equiv \frac{p(x|\theta_0)}{p(x|\theta_1)}, \quad (4)$$

an otherwise intractable quantity. This motivates a family of new techniques for likelihood-free inference in which the the joint likelihood ratio and joint score are used as training data for neural networks. These networks serve as surrogate models for the intractable likelihood or likelihood ratio. Experiments showed these new methods to be more sample-efficient than previously established neural density and neural density ratio estimation techniques. The authors of Refs. [24–26] coined the term “mining gold” for the process of extracting the joint likelihood ratio and joint score from the simulator – while the augmented data require some effort to extract, they are extremely valuable.

While the loss functionals originally proposed in Refs. [24–26] have the correct minima, they are not necessarily the most sample efficient. In particular, the proposed mean squared error (MSE) losses are often dominated by few samples with large joint likelihood ratios. Here we extend and improve that original work with two new algorithms for likelihood-free inference. The key improvement are new loss functions, which use an improved estimator for the cross entropy based on the joint likelihood ratio and joint score.

After introducing these new algorithm in Sec. 2, we show its performance in a problem from particle physics in Sec. 3, before giving our conclusions in Sec. 4.

## 2 Cross-entropy estimation with augmented data

Consider the problem of estimating the likelihood ratio  $r(x|\theta_0, \theta_1)$  based on samples  $(x_i, z_i) \sim p(x, z|\theta_0)$ , labeled with  $y_i = 0$ ; samples  $(x_i, z_i) \sim p(x, z|\theta_1)$ , labeled  $y_i = 1$ ; and the joint likelihood ratio  $r(x_i, z_i|\theta_0, \theta_1)$  and joint score  $t(x_i, z_i|\theta_0)$ .

The familiar binary cross-entropy loss functional is defined as

$$L[\hat{s}(x)] = - \int dx \left[ p(x|y=1) \log(\hat{s}(x)) + p(x|y=0) \log(1 - \hat{s}(x)) \right]. \quad (5)$$

For balanced samples ( $p(\theta_0) = p(\theta_1) = 1/2$ ) we have

$$p(x, z) = \frac{p(x, z|\theta_0) + p(x, z|\theta_1)}{2} \quad (6)$$

$$s(x, z|\theta_0, \theta_1) = p(y=1|x, z) = \frac{1}{r(x, z|\theta_0, \theta_1) + 1} = \frac{p(x, z|\theta_1)}{p(x, z|\theta_0) + p(x, z|\theta_1)}, \quad (7)$$

which allows us to rewrite Eq. 5 as

$$L[\hat{s}(x)] = - \int dx dz p(x, z) \left[ s(x, z|\theta_0, \theta_1) \log(\hat{s}(x)) + (1 - s(x, z|\theta_0, \theta_1)) \log(1 - \hat{s}(x)) \right]. \quad (8)$$

It is straightforward to show that this loss functional is minimized by

$$s^*(x) \equiv \arg \min_{\hat{s}(x)} L[\hat{s}(x)] = s(x|\theta_0, \theta_1) = p(y=1|x) = \frac{1}{r(x|\theta_0, \theta_1) + 1} = \frac{p(x|\theta_1)}{p(x|\theta_0) + p(x|\theta_1)}. \quad (9)$$

To use the cross entropy to train a surrogate model with a finite number of samples, we need a tractable estimator for the cross entropy. The standard estimator, as used for instance in the CARL inference method [9], is given by

$$L_{\text{CARL}}[\hat{s}(x)] = - \frac{1}{N} \sum_{(x_i, y_i)} \left[ y_i \log(\hat{s}(x_i)) + (1 - y_i) \log(1 - \hat{s}(x_i)) \right]. \quad (10)$$

The  $y_i \in \{0, 1\}$  act as an unbiased, but high-variance estimator of  $s(x_i, z_i|\theta_0, \theta_1)$ . In the limit of infinite samples, this estimator therefore has the correct minimum of Eq. (9), but for finite sample sizes it may suffer from high variance.

With the availability of the joint likelihood ratio  $r(x_i, z_i|\theta_0, \theta_1)$  from the simulator, the  $s(x_i, z_i|\theta_0, \theta_1)$  are tractable and we can define the alternative estimator

$$L_{\text{ALICE}}[\hat{s}(x)] = -\frac{1}{N} \sum_{(x_i, z_i) \sim p(x_i, z_i)} \left[ s(x_i, z_i|\theta_0, \theta_1) \log(\hat{s}(x_i)) + (1-s(x_i, z_i|\theta_0, \theta_1)) \log(1-\hat{s}(x_i)) \right]. \quad (11)$$

By using the exact  $s(x, z|\theta_0, \theta_1)$  rather than the  $y_i \in \{0, 1\}$ , the samples drawn according to  $y = 0$  also provide information about the second  $y = 1$  term in the loss function, and vice versa. By minimizing the loss function we get an estimator  $\hat{s}(x)$  and thus a likelihood ratio estimator

$$\hat{r}(x|\theta_0, \theta_1) = \frac{1 - \hat{s}(x)}{\hat{s}(x)}. \quad (12)$$

This defines the ALICE inference method<sup>1</sup>, which consists of mining the joint likelihood ratio from the simulator, training a neural network on the improved cross-entropy estimator in Eq. (11), and using this surrogate model for statistical inference on  $\theta$ .

It is to be expected that a likelihood ratio estimator based on the ALICE estimator for the cross-entropy should outperform the CARL method, which is based on the standard cross-entropy estimator in Eq. (10). The more interesting question is how it stacks up against the ROLR technique introduced in Refs. [24–26], in which the loss function

$$L_{\text{ROLR}}[\hat{r}(x)] = \frac{1}{N} \sum_{(x_i, y_i, z_i)} \left[ y_i |r(x_i, z_i|\theta_0, \theta_1) - \hat{r}(x|\theta_0, \theta_1)|^2 + (1 - y_i) \left| \frac{1}{r(x_i, z_i|\theta_0, \theta_1)} - \frac{1}{\hat{r}(x|\theta_0, \theta_1)} \right|^2 \right] \quad (13)$$

is minimized. In the limit of infinite samples it is minimized by  $r(x|\theta_0, \theta_1)$ . But here each event only contributes to either the squared error on  $r$  or on  $1/r$  term, which might lead to a higher variance.

In analogy to the CASCAL and RASCAL methods of Refs. [24–26], we can define an additional inference method which uses the joint score, i. e. an additional piece of information that describes the local (tangential) behavior of the likelihood function. If a parameterized likelihood ratio estimator is implemented with a differentiable architecture such as a neural network, we can calculate the gradient of the output  $\hat{s}(x|\theta_0, \theta_1)$  with respect to  $\theta_0$  and similarly calculate the corresponding score

$$\hat{t}(x|\theta_0, \theta_1) = \nabla_{\theta} \log \hat{r}(x|\theta_0, \theta_1) = \nabla_{\theta} \log \left( \frac{1 - \hat{s}(x_i|\theta, \theta_1)}{\hat{s}(x_i|\theta, \theta_1)} \right) \quad (14)$$

of the  $\hat{r}$  estimator. For a perfect  $\hat{r}$  (or equivalently  $\hat{s}$ ) estimator, this corresponding score  $\hat{t}$  will also minimize the squared error loss with respect to the joint score  $t(x, z|\theta_0, \theta_1)$ , which can be extracted from the simulator [24–26]. Turning this argument around, we can use the joint score to guide the training of the estimator. This is the idea behind the ALICES<sup>2</sup> technique, which is based on the loss function

$$L_{\text{ALICES}}[\hat{s}(x|\theta_0, \theta_1)] = -\frac{1}{N} \sum_{(x_i, z_i) \sim p(x_i, z_i)} \left[ s(x_i, z_i|\theta_0, \theta_1) \log(\hat{s}(x_i)) + (1 - s(x_i, z_i|\theta_0, \theta_1)) \log(1 - \hat{s}(x_i)) + \alpha (1 - y_i) \left| t(x_i, z_i|\theta_0, \theta_1) - \nabla_{\theta} \log \left( \frac{1 - \hat{s}(x_i|\theta, \theta_1)}{\hat{s}(x_i|\theta, \theta_1)} \right) \right|_{\theta_0}^2 \right]. \quad (15)$$

The factor  $(1 - y_i)$  is necessary to guarantee the correct minimum of the squared error on the score. The hyper-parameter  $\alpha$  weights the two terms in the loss function. This loss is the natural extension of the the CASCAL loss function, but we expect it to reduce the variance compared to the CASCAL approach for finite sample size. An interesting question is how it performs compared to the RASCAL approach, which similarly augments the ROLR loss in Eq. (13) with the score term.

<sup>1</sup>Approximate likelihood with improved cross-entropy estimator

<sup>2</sup>Approximate likelihood with improved cross-entropy estimator and score

Strategy	Expected MSE		
	$10^4$ training samples	$10^5$ training samples	$10^7$ training samples
Histogram			0.0561
CARL	<b>0.1743</b>	<b>0.1672</b>	<b>0.0124</b>
ROLR	0.1345	0.0396	0.0032
CASCAL	0.1715	0.1652	0.0008
RASCAL	0.0449	0.0100	0.0009
ALICE	0.0510	<b>0.0076</b>	<b>0.0004</b>
ALICES	<b>0.0339</b>	0.0111	0.0013
SALLY	<b>0.0261</b>	<b>0.0146</b>	<b>0.0132</b>
SALLINO	0.0319	0.0227	0.0213

Table 1: Fidelity of different strategies in the idealized scenario, using training sets of various sizes. We use the expected mean squared error as defined in Ref. [25] as a performance metric. The new methods, ALICE and ALICES, outperform RASCAL, CASCAL, and ROLR. The SALLY and SALLINO methods are very sample efficient, but make approximations that limit their asymptotic performance.

### 3 Experiments

We experiment with the new methods in the particle physics problem introduced in Refs. [24, 25]. In this real-world problem, the outcome of proton-proton collisions is characterized by 42 observables, from which likelihood ratios and confidence limits on two model parameters are derived. We first consider an idealized setting neglecting the detector response where the likelihood function is tractable, which provides us with ground truth that can be used to evaluate the performance of the algorithms. For a detailed description of the setup, see Ref. [25].

In addition to the CARL, ROLR, CASCAL, and RASCAL techniques described above, we also compare to the SALLY and SALLINO methods. SALLY and SALLINO approximate a statistical model that is accurate in the neighborhood of  $\theta = (0, 0)^T$ . The methods are very sample efficient, but make approximations that limit their asymptotic performance.

Except for the new loss functions, we used the same architectures and hyper-parameters as in Ref. [25]. In particular, we use fully connected networks with five hidden layers, 100 units each, and tanh activation functions for both approaches. For ALICES we use  $\alpha = 5$ , which was found to give a good performance for the closely related CASCAL method [25].

Table 1 and Fig. 1 show the quality of the likelihood ratio estimate based on various sized training samples for the new methods and compares them to the inference techniques presented in Ref. [25]. As a performance metric we use an expected mean squared error on the log likelihood ratio, as defined in Ref. [25].

Unsurprisingly, the ALICE and ROLR methods clearly outperform CARL, which does not have access to the joint likelihood ratio. More significantly, we find that ALICE outperforms ROLR, which does have access to the joint likelihood ratio. We conjecture that this improvement can be attributed to the lower variance of the cross-entropy compared to the squared error. More surprisingly, the ALICE method also outperforms the RASCAL method for larger training sample sizes ( $\geq 10^5$ ), even though ALICE does not have access to the joint score.

For smaller training sample sizes ( $\leq 10^5$ ) the ALICES method outperforms the ALICE method, which is not surprising given the additional information available during training. For larger training sample sizes ( $\geq 10^5$ ), the variance of the score actually deteriorates the performance of ALICES compared to ALICE. We did not perform hyper-parameter tuning for  $\alpha$  as a function of the training sample size, which should ensure that ALICES performs at least as well as ALICE. We leave a systematic tuning of the  $\alpha$  parameter and an analysis of sources of variance in this approach for future work.

Figure 2 shows expected exclusion contours at different confidence levels on the two parameters, assuming 36 observed events distributed according to  $\theta = (0, 0)^T$ . The methods are trained on the full training samples of  $10^7$  samples. The left panel shows contours constructed based on asymptotic

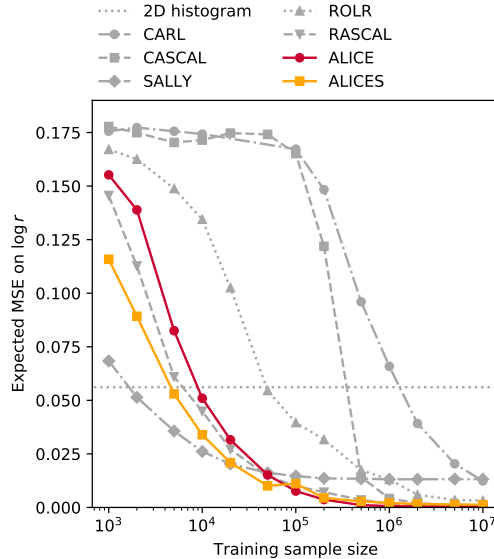


Figure 1: Estimator fidelity in the idealized scenario as a function of the training sample size. As a metric we use the expected mean squared error on the log likelihood ratio, see Ref. [25]. The new methods are more sample efficient than the similar ROLR and RASCAL techniques.

properties of the profile likelihood ratio test statistic. While methods such as RASCAL are generally very accurate, with this construction they can sometimes lead to overly optimistic exclusion contours, visible as tighter bounds than the “truth” contour. We find that switching to ALICE and ALICES reduces this issue, but does not entirely solve it.

The right panel of Fig. 2 shows exclusion contours based on the frequentist confidence intervals calibrated with toy experiments. This Neyman construction guarantees coverage: while the limits from any approach may be worse than the optimal limits, they will never be overly optimistic. As test statistics we use the likelihood ratio with respect to the  $\theta = (0, 0)^T$ , which explains why the contours are generally stronger than in the left panel. We find that both ALICE and ALICES, like RASCAL and CASCAL of Refs. [24–26], lead to limits that are virtually indistinguishable from the ideal limits based on the true likelihood ratio.

Finally, in Fig. 3 we show similar expected exclusion contours, but in a more realistic setup in which the parton shower and detector effects are described with approximate smearing functions, which makes the true likelihood intractable. In this situation, we cannot compare the likelihood ratio estimators to the ground truth. Instead, we show the expected contours based on the Neyman construction, similar to the right panel of Fig. 2. In the left panel we show results for limited training samples of only  $10^5$  events. In this setup, ALICES allows for strong limits, comparable to RASCAL and slightly better than for ALICE. The right panel demonstrates that with the full training sample the results of RASCAL, ALICE, and ALICES are indistinguishable.

## 4 Conclusions

In this work, we have extended recently developed inference techniques for the setting in which the likelihood is only implicitly defined through a stochastic generative model or simulator. By exploiting the joint likelihood ratio that can be extracted from the simulator, we introduced an improved cross-entropy estimator. This improved cross-entropy estimator is used to define two new likelihood-free inference techniques: ALICE and ALICES.

Our experiments comparing ALICE and ALICES with the other recently developed techniques indicate that they are significantly more sample efficient than ROLR, CASCAL, and RASCAL techniques. We attribute this to the lower variance of the improved cross-entropy estimator. For smaller training sample sizes, there are still advantages to the SALLY and SALLINO techniques.

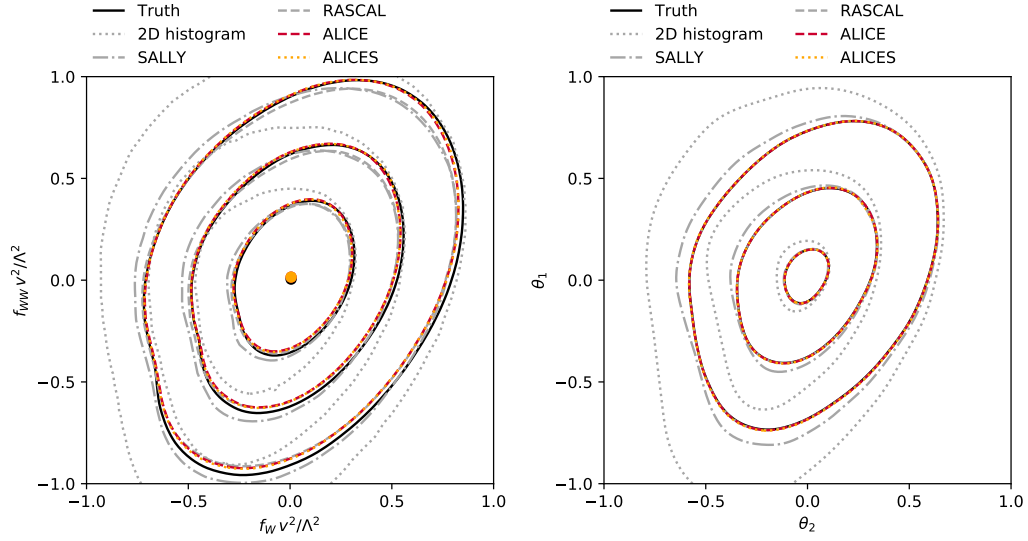


Figure 2: Expected exclusion limits on the model parameters in the idealized scenario for different inference methods. We assume 36 events distributed according to  $\theta = (0, 0)^T$ . All estimators are trained on a large data set with  $10^7$  samples. Left: construction of exclusion limits based on asymptotic properties of the likelihood ratio. With this method, inefficient estimators can predict overly optimistic exclusion limits, as can be seen for instance for the RASCAL method. The new ALICE and ALICES approaches are less prone to this issue. Right: construction of exclusion limits calibrated with toy experiments (i. e. the Neyman construction). In this approach, the intervals will always cover, but might not be optimal. We find an excellent performance of the ALICE and ALICES methods, virtually indistinguishable from the RASCAL method and the true likelihood ratio.

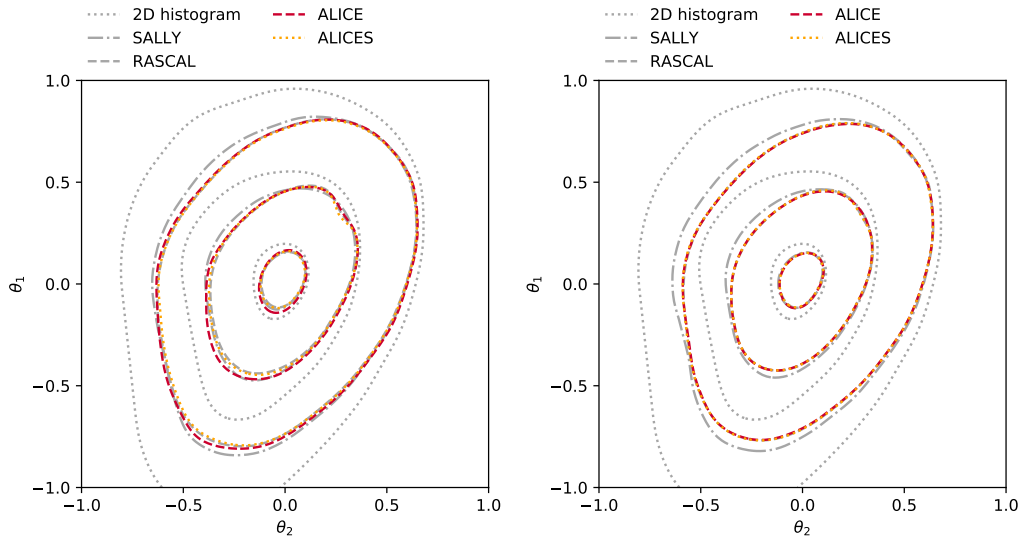


Figure 3: Expected exclusion limits on the model parameters in the scenario with detector effects for different inference methods. We construct the contours with the Neyman construction, which guarantees the confidence intervals will cover. The intervals are based on 36 events distributed according to  $\theta = (0, 0)^T$ . The estimators are trained on data sets with  $10^5$  (left) or  $10^7$  (right) samples. The ALICES method leads to strong limits, comparable to the RASCAL technique.

We note that it is possible to use a hybrid of the traditional cross-entropy of Eq. 5 and the improved cross-entropy Eq. 11. This would be useful in situations where one may not have access to the joint ratio for practical reasons or because some training samples come from real data instead of a simulation. Furthermore, we note that the improved cross-entropy estimator of ALICE and ALICES can be extended from the binary setting to one where samples are generated from multiple parameters points if the joint likelihood ratio for all pairs is available. These joint likelihood ratios provide the necessary ingredient for importance sampling beyond the binary setting considered here.

The ubiquity of simulators and other implicit models indicates there is enormous potential for likelihood-free inference techniques. The use of augmented data improves the sample efficiency of these techniques significantly, and these results motivate further study of variance reduction techniques that leverage this augmented data.

## Acknowledgments

JB, KC, and GL are grateful for the support of the Moore-Sloan data science environment at NYU. KC and GL were supported through the NSF grants ACI-1450310 and PHY-1505463. JP was partially supported by the Scientific and Technological Center of Valparaíso (CCTVal) under Fondecyt grant BASAL FB0821. This work was supported in part through the NYU IT High Performance Computing resources, services, and staff expertise.

## References

- [1] D. B. Rubin: ‘Bayesianly justifiable and relevant frequency calculations for the applied statistician’. *Ann. Statist.* 12 (4), p. 1151, 1984. URL <https://doi.org/10.1214/aos/1176346785>.
- [2] M. A. Beaumont, W. Zhang, and D. J. Balding: ‘Approximate bayesian computation in population genetics’. *Genetics* 162 (4), p. 2025, 2002.
- [3] J. Alsing, B. Wandelt, and S. Feeney: ‘Massive optimal data compression and density estimation for scalable, likelihood-free inference in cosmology’, 2018. arXiv:1801.01497.
- [4] T. Charnock, G. Lavaux, and B. D. Wandelt: ‘Automatic physical inference with information maximizing neural networks’. *Phys. Rev. D* 97 (8), p. 083004, 2018. arXiv:1802.03537.
- [5] T. Kanamori, S. Hido, and M. Sugiyama: ‘A least-squares approach to direct importance estimation’. *Journal of Machine Learning Research* 10 (Jul), p. 1391, 2009.
- [6] Y. Fan, D. J. Nott, and S. A. Sisson: ‘Approximate Bayesian Computation via Regression Density Estimation’. *ArXiv e-prints*, 2012. arXiv:1212.1479.
- [7] L. Dinh, D. Krueger, and Y. Bengio: ‘NICE: Non-linear Independent Components Estimation’. *ArXiv e-prints*, 2014. arXiv:1410.8516.
- [8] D. Jimenez Rezende and S. Mohamed: ‘Variational Inference with Normalizing Flows’. *ArXiv e-prints*, 2015. arXiv:1505.05770.
- [9] K. Cranmer, J. Pavez, and G. Louppe: ‘Approximating Likelihood Ratios with Calibrated Discriminative Classifiers’, 2015. arXiv:1506.02169.
- [10] K. Cranmer and G. Louppe: ‘Unifying generative models and exact likelihood-free inference with conditional bijections’. *J. Brief Ideas*, 2016.
- [11] L. Dinh, J. Sohl-Dickstein, and S. Bengio: ‘Density estimation using Real NVP’. *ArXiv e-prints*, 2016. arXiv:1605.08803.
- [12] G. Papamakarios and I. Murray: ‘Fast  $\epsilon$ -free inference of simulation models with bayesian conditional density estimation’. In ‘Advances in Neural Information Processing Systems’, p. 1028–1036, 2016.
- [13] B. Paige and F. Wood: ‘Inference Networks for Sequential Monte Carlo in Graphical Models’. *ArXiv e-prints*, 2016. arXiv:1602.06701.

- [14] R. Dutta, J. Corander, S. Kaski, and M. U. Gutmann: ‘Likelihood-free inference by ratio estimation’. ArXiv e-prints , 2016. arXiv:1611.10242.
- [15] B. Uria, M.-A. Côté, K. Gregor, I. Murray, and H. Larochelle: ‘Neural Autoregressive Distribution Estimation’. ArXiv e-prints , 2016. arXiv:1605.02226.
- [16] A. van den Oord, S. Dieleman, H. Zen, et al.: ‘WaveNet: A Generative Model for Raw Audio’. ArXiv e-prints , 2016. arXiv:1609.03499.
- [17] A. van den Oord, N. Kalchbrenner, O. Vinyals, L. Espeholt, A. Graves, and K. Kavukcuoglu: ‘Conditional Image Generation with PixelCNN Decoders’. ArXiv e-prints , 2016. arXiv:1606.05328.
- [18] A. van den Oord, N. Kalchbrenner, and K. Kavukcuoglu: ‘Pixel Recurrent Neural Networks’. ArXiv e-prints , 2016. arXiv:1601.06759.
- [19] M. U. Gutmann, R. Dutta, S. Kaski, and J. Corander: ‘Likelihood-free inference via classification’. *Statistics and Computing* p. 1–15, 2017.
- [20] D. Tran, R. Ranganath, and D. M. Blei: ‘Hierarchical Implicit Models and Likelihood-Free Variational Inference’. ArXiv e-prints , 2017. arXiv:1702.08896.
- [21] G. Louppe and K. Cranmer: ‘Adversarial Variational Optimization of Non-Differentiable Simulators’. ArXiv e-prints , 2017. arXiv:1707.07113.
- [22] G. Papamakarios, T. Pavlakou, and I. Murray: ‘Masked Autoregressive Flow for Density Estimation’. ArXiv e-prints , 2017. arXiv:1705.07057.
- [23] G. Papamakarios, D. C. Sterratt, and I. Murray: ‘Sequential Neural Likelihood: Fast Likelihood-free Inference with Autoregressive Flows’. ArXiv e-prints , 2018. arXiv:1805.07226.
- [24] J. Brehmer, K. Cranmer, G. Louppe, and J. Pavez: ‘Constraining Effective Field Theories with Machine Learning’ , 2018. arXiv:1805.00013.
- [25] J. Brehmer, K. Cranmer, G. Louppe, and J. Pavez: ‘A Guide to Constraining Effective Field Theories with Machine Learning’ , 2018. arXiv:1805.00020.
- [26] J. Brehmer, G. Louppe, J. Pavez, and K. Cranmer: ‘Mining gold from implicit models to improve likelihood-free inference’ , 2018. arXiv:1805.12244.