# Stochastic Gradient Bayesian Optimal Experimental Designs for Simulation-based Inference

**Vincent D. Zaballa** [1]  **Elliot E. Hui** [1]

## Abstract

Simulation-based inference (SBI) methods tackle complex scientific models with challenging inverse problems. However, SBI models often face a significant hurdle due to their non-differentiable nature, which hampers the use of gradient-based optimization techniques. Bayesian Optimal Experimental Design (BOED) is a powerful approach that aims to make the most efficient use of experimental resources for improved inferences. While stochastic gradient BOED methods have shown promising results in high-dimensional design problems, they have mostly neglected the integration of BOED with SBI due to the difficult non-differentiable property of many SBI simulators. In this work, we establish a crucial connection between ratio-based SBI inference algorithms and stochastic gradient-based variational inference by leveraging mutual information bounds. This connection allows us to extend BOED to SBI applications, enabling the simultaneous optimization of experimental designs and amortized inference functions. We demonstrate our approach on a simple linear model and offer implementation details for practitioners.

## 1. Introduction

Many scientific models are defined by a simulator that defines an output $y$ determined by the inputs, or designs, to a system, $\xi$, and parameters that define how the model transforms the inputs to outputs, $\theta$. Inferring a distribution of parameters given data $p(\theta|y,\xi)$ is of central importance in Bayesian statistics and can be seen as a form of solving the inverse problem for a given simulator(Lindley, 1972). In SBI, a simulator forms an implicit probability distribu-

tion known as the likelihood $p(y|\theta,\xi)$ that is used with the prior of the model parameters $p(\theta)$ to estimate the posterior probability of the model parameters given the observed data, $p(\theta|y,\xi)$ (Cranmer et al., 2020). Recent SBI methods have use deep learning-based models to infer either the intractable likelihood or posterior using density estimators for both, or classifiers to estimate the likelihood-to-evidence ratio, $\frac{p(\theta|y,\xi)}{p(\theta|\xi)} = \frac{p(y|\theta,\xi)}{p(y|\xi)} = \frac{p(y,\theta|\xi)}{p(\theta)p(y|\xi)}$.

However, inferring the likelihood, posterior, or ratio is a computationally expensive process that depends on observed data $y_o$, to compute. Recent work questioned the validity of this expensive computational process used in SBI if using the wrong simulator for the true data generating process (Cannon et al.). Naïve conclusions can be made if using the wrong model of the underlying scientific phenomenon, or the model is not close enough to the real data generating process, which motivates the use of optimal experimental designs in SBI methods.

Bayesian optimal experimental design (BOED) has shown promise as a way to optimize experiments given a model, the simulator, and priors of the parameters of interest. BOED works by determining the information gain of a proposed experimental design, $\xi$, on the parameters of the model of interest

$$IG(y,\xi) = H[p(\theta)] - H[p(\theta|y,\xi)]. \quad (1)$$

The information gain can only be evaluated after an experiment but another quantity, the Expected Information Gain (EIG), $I(\xi)$, can be used as a proxy for the information gained in an experiment

$$I(\xi) \triangleq \mathbb{E}_{p(y|\xi)}\left[H[p(\theta)] - H[p(\theta|y,\xi)]\right], \quad (2)$$

The intuition behind this process is we must ask ourselves, which experimental design and outcome would be most surprising given what we assume, or know, about the system when conducting the experiment. This can be rewritten into the form of calculating the mutual information between the observed data and unknown parameters

$$I(\xi) = \text{MI}_\xi(\theta;y) = \mathbb{E}_{p(\theta)p(y|\theta,\xi)}\left[\log\frac{p(y|\theta,\xi)}{p(y|\xi)}\right]. \quad (3)$$

Early BOED work focused on estimating the mutual information then using that estimate as the surrogate function in

[1]Department of Biomedical Engineering, University of California, Irvine, Irvine, CA, USA. Correspondence to: Vincent Zaballa <vzaballa@uci.edu>, Elliot Hui <eehui@uci.edu>.

an outer optimizer, such as Bayesian optimization (Foster et al., 2019b; Kleinegesse & Gutmann, 2019). This double loop of optimization was inefficient and lead to development of methods to simultaneously optimize the design and mutual information in a single optimization process. However, this unified optimization depended on an unnormalized likelihood and posterior approximation (Foster et al., 2019a) or an implicit likelihood with a simulator that has a differentiable functional form (Kleinegesse & Gutmann, 2021).

We present a method to simultaneously optimize designs and the mutual information for the remaining set of models, implicit likelihoods without a differentiable simulator, which are typically used in the SBI literature. We additionally make a link to how we can use a generative model in Contrastive Precitive Coding. We show:

- A differentiable objective for simultaneously optimizing the mutual information and likelihood for SBI-based models.

- A connection between Likelihood-Free based methods for BOED and contrastive ratio estimation (CRE) methods for SBI models.

- Experimental validation of the unified objective on a simple linear model.

## 2. Background

Previous work in SBI methods have focused on improving methods based on given, observed, data $y_o$, (Papamakarios & Murray, 2016; Papamakarios et al., 2018; Durkan et al., 2020; Greenberg et al., 2019) whereas BOED has focused on determining an optimal design $\xi^*$, based on various bounds of MI between $y$ and $\theta$. While these aims seem to be unrelated, we will show how they can be performed simultaneously for SBI methods that rely on potentially stochastic simulators that act as black-box functions.

### 2.1. Simulation-Based Inference

In many scientific disciplines, it is desirable to infer a distribution of parameters $\theta$, of a potentially stochastic model, or simulator, given observations, $y_o$. The closed-box simulator may depend on random numbers $z$, such as in stochastic differential equations, and previous experimental designs $\xi$, such that the simulator takes the form $y = g(\theta, \xi, z)$. When a likelihood is not available, Approximate Bayesian Computation (ABC) methods can be used, (Sisson et al., 2018) which aim to infer the likelihood of parameters of the simulator that are within an $\epsilon$ ball, $B_\epsilon(y)$, of the observed data $y := y_o$, resulting in the likelihood $p(\|y - y_o\| < \epsilon|\theta)$. However, recent SBI methods have outperform ABC methods in inference tasks (Lueckmann et al., 2021). By using a simulator to simulate the joint data distribution $(\theta, y) \sim p(y|\theta)$,

drawn from a prior $\theta \sim p(\theta)$, we can obtain an amortized likelihood $p_\phi(y|\theta)$ or posterior $p_\phi(\theta|y)$ by training a neural density estimator, such as a normalizing flow, with parameters $\phi$, or estimate the likelihood-to-evidence ratio $\exp f_\phi(\theta, y) \approx \frac{p(y|\theta)}{p(y)}$, by training a classifier to distinguish parameters used to simulate an observed values, $y$. Different SBI methods can be used in inference for downstream applications depending on the desiderata of the inference task. For example, one might use an amortized posterior approximation if there are many different data samples to evaluate, whereas an ensemble of ratios (Hermans et al.) has been shown to perform more robustly on Simulation-Based Calibration (SBC) tests (Talts et al.) at the cost of increased computational complexity.

There are many SBI methods proposed for approximating the likelihood, posterior, or ratio. We review the relevant ones to our method here. See (Lueckmann et al., 2021) for a more thorough review and benchmark of SBI methods.

**Neural Likelihood Estimation** We can use data from the joint distribution to train a conditional neural density-based likelihood function. If we take a dataset of samples $\{y_n, \theta_n\}_{1:N}$ obtained from a simulator as previously described, we can train a conditional density estimator $p_\phi(y|\theta)$ to model the likelihood by maximizing the total log likelihood of $\sum_n \log p_\phi(y_n|\theta_n)$, which is approximately equivalent to minimizing the loss

$$\mathcal{L}(\phi) = \mathbb{E}_{p(\boldsymbol{\theta})}(D_{\mathrm{KL}}(p(y|\theta)\|p_\phi(y|\theta))) + \mathrm{const}, \quad (4)$$

where the Kullback-Leibler divergence is minimized when $p_\phi(y|\theta)$ approaches $p(y|\theta)$. SBI methods would then condition this likelihood on observed data, $y_o$, and refine the likelihood estimate by resetting the prior to become the new posterior samples via Markov Chain Monte Carlo (MCMC) sampling of the approximate likelihood $p(\theta) := p(\theta|y_o) \propto p_\phi(y_o|\theta)p(\theta)$ and training a new neural density estimator of the likelihood (Papamakarios et al., 2018; Lueckmann et al., 2018). This is Sequential Neural Likelihood (SNL) which we forego as we focus on the preliminary step of optimizing an experimental design without $y_o$.

### 2.2. Bayesian Optimal Experimental Design

Following from equation 3, (Foster et al., 2019a) proposed the prior contrastive estimation (PCE) lower bound of the MI

$$I_{PCE}(\xi, L) \triangleq \mathbb{E}\left[\log \frac{p(y|\theta_0, \xi)}{\frac{1}{L+1}\sum_{\ell=0}^{L} p(y|\theta_\ell, \xi)}\right], \quad (5)$$

where the expectation is over $p(\theta_0)p(y|\theta_0, \xi)p(\theta_{1:L})$ and $\xi$ is the proposed design, $\theta_0$ is the original parameter that generated data $y$, and $L$ is the number of contrastive samples. The PCE bound is appropriate in BOED when the

prior and posterior are similar enough that $p(\theta)$ is a suitable proposal distribution for $p(y|\xi)$. This bound has low variance but is upper-bounded by $\log L$, potentially leading to large bias but still demonstrated adequate performance on various benchmarks. Unfortunately, this bound requires a tractable likelihood function, which is not available in SBI applications.

## 3. SBI-based BOED

### 3.1. Likelihood Free PCE

We take inspiration from previous SBI and BOED methods to allow optimization of designs with respect to closed-box simulators that are modeled using normalizing flows. We start by noting how the loss function of contrastive ratio estimation (CRE) (Durkan et al., 2020) lower bounds PCE

$$
\log \frac{\exp(f_\phi(\theta, y))}{\sum_{\ell=0}^{L} \exp(f_\phi(\theta_\ell, y))} \leq \log \frac{\exp(f_\phi(\theta, y))}{\frac{1}{1+L} \sum_{\ell=0}^{L} \exp(f_\phi(\theta_\ell, y))}
$$
$$
= \log \frac{p_\phi(y|\theta_0, \xi)}{\frac{1}{1+L} \sum_{\ell=0}^{L} p_\phi(y|\theta_l, \xi)},
$$

where $L$ is the number of contrastive samples and $f_\phi$ is a discriminative classifier, which holds for a single batch of data and constant experimental design, i.e. when $\xi$ is constant. We exchange an explicit likelihood in PCE with a neural density estimator to create Likelihood-Free PCE (LF-PCE). We now have a MI lower bound

$$
I(\xi, \phi, L) \geq \mathbb{E}\left[\log \frac{p_\phi(y|\theta_0, \xi)}{\frac{1}{1+L} \sum_{\ell=0}^{L} p_\phi(y|\theta_l, \xi)}\right], \quad (6)
$$

where the expectation is over $p(\theta_0)p(y|\theta_0, \xi)p(\theta_{1:L})$. We now can simultaneously optimizes designs and parameters of a neural density estimator. If we are to use a normalizing flow as $\exp f_\phi(y, \theta, \xi) = p_\phi(y|\theta, \xi)$, then the PCE lower bound of the MI holds since the distribution is normalized as normalizing flows are bounded functions (Papamakarios et al., 2019). We note that this can be an unstable objective as the data distribution of the flow will change as experimental designs change. However, the result is that it returns an amortized likelihood that can be evaluated on observed experimental data to return a posterior density or used in downstream inference algorithms, such as SNL. Finally, using a normalizing flow allows us to take gradients with respect to designs $\xi$, which we derive in Appendix A.

**Practical implementation of LF-PCE loss** For LF-PCE training, stability of the density estimator is a challenge when optimizing the MI lower bound. To address this, we added a regularization term, $\lambda$, to both loss functions to help stabilize the training of the density estimator during design

optimization

$$
\mathbb{E}\left[\log \frac{p_\phi(y|\theta_0, \xi)}{\frac{1}{1+L} \sum_{\ell=0}^{L} p_\phi(y|\theta_l, \xi)} + \lambda \cdot \log p_\phi(y|\theta_0, \xi)\right], \quad (7)
$$

where the expectation is over $p(\theta_0)p(y|\theta_0, \xi)p(\theta_{1:L})$.

### 3.2. Connection to Generative MI Estimation

The mutual information bound proposed by (Foster et al., 2019a) for PCE is similar to Contrastive Predictive Coding (CPC) (Poole et al., 2019; Oord et al., 2018), but where a generative model replaces a discriminative one and the random variable X corresponds to observed data and random variable Y to the prior distribution. In our formulation the bound of the MI depends on both the amount of training $tr \rightarrow \infty$ and number of contrastive samples $L \rightarrow \infty$ to approach the true MI. The generative approach to CPC can be simplified as

$$
I_{PCE}(\phi) := \mathbb{E}_P[\log p_\phi(x|y) - \log p_\phi(x)], \quad (8)
$$

where $P$ is a random variable representing the joint distribution we obtain from our simulators $(x, y) \sim p(x|y)p(y)$ and $p_\phi(x)$ implicitly depends on the number of contrastive samples $L$ to approximate the marginal likelihood.

## 4. Experimental Evaluation

### 4.1. Noisy Linear Model

We follow (Kleinegesse & Gutmann, 2020) and evaluate optimal designs on a classic noisy linear model where a response variable $y$ has a linear relationship with experimental designs $\xi$, which is determined by values of the model parameters $\boldsymbol{\theta} = [\theta_0, \theta_1]$, which model the offset and gradient. We would like to optimize the value of $D$ measurements to estimate the posterior of $\boldsymbol{\theta}$, and so create a design vector $\boldsymbol{\xi} = [\xi_1, \ldots, \xi_D]^{\mathsf{T}}$. Each design, $\xi_i$ returns a measurement $y_i$, which results in the data vector $\mathbf{y} = [y_1, \ldots, y_D]^{\mathsf{T}}$. We assume non-Gaussian noise sources, otherwise evaluating the posterior and MI would be trivial. We use a Gaussian noise source $\mathcal{N}(\epsilon; 0, 1)$ and Gamma noise source $\Gamma(\nu; 2, 2)$. The model is then

$$
y = \theta_0 \mathbf{1} + \theta_1 * \boldsymbol{\xi} + \boldsymbol{\epsilon} + \boldsymbol{\nu}, \quad (9)
$$

where $\boldsymbol{\epsilon} = [\epsilon_1, \ldots, \epsilon_D]^{\mathsf{T}}$ and $\boldsymbol{\nu} = [\nu_1, \ldots, \nu_D]^{\mathsf{T}}$ are i.i.d. samples. We evaluate LF-PCE on this model and examine how changing the $\lambda$ regularization parameter in (7) influences the resulting mutual information bound and design quality for both models.

For each design dimension, D, we randomly initialize designs $\xi \in [-10, 10]$. For LF-PCE, we chose $N = 10$, the number of non-contrastive samples $y \sim p(y|\xi, \theta_0)$, and

*Figure 1.* Comparison of the EIG across design dimensions, type of BOED, and $\lambda$ regularization for the noisy linear model examining the moving average over N=10 samples. For the single design dimension, LF-PCE with no $\lambda$ regularization outperforms in estimating a lower bound of the MI, which can translate to more informative experimental designs. In the higher-dimension design cases, LF-PCE increases its EIG with more designs, which is expected, but sees diminishing returns when expanding from 10D to 100D design evaluations. In the 100-dimensional design case, we see the benefit of using $\lambda$ regularization to stabilize the training of a neural density estimator in high-dimensional input space at the cost of slightly lower EIG.

$M = 50$ contrastive samples for all experiments. For the neural spline flow, we chose 5 bijector layers, each with 4 bins, and 4 resnet multilayer perceptrons, each with 128 dimensions, for the neural network-based conditional networks. For both the neural density estimator's parameters $\phi$, and the designs $\xi$, we use the Adam optimizer (Kingma & Ba, 2014) with $\beta_1 = 0.9$ and $\beta_2 = 0.99$, with learning rate $\alpha = 1e^{-3}$ for the neural density estimator and $\alpha = 1e^{-2}$ for design optimization.

Examining the graph of the mutual information in Figure 1, we see that LF-PCE lower bound steadily increases for all values of lambda; however, the stability of the optimization of the generative model's parameters diverges in higher design dimensions whenever $\lambda = 0$. We see a general trend between exploration and exploitation in changing values of $\lambda$, where higher $\lambda$ values lead to lower MI lower-bound estimates and potentially more homogenous designs.

Using LF-PCE we obtain an amortized neural density estimator of the likelihood that is able to perform inference on observed data evaluated at the optimal design. For example, $p(\theta|y_o, \xi^*) \propto p_\phi(y_o|\theta, \xi^*)p(\theta)$ by MCMC sampling. We evaluate the posterior densities after optimizing on the LF-PCE lower bound in Appedix B and can see the mean and interquartile range in Table 1. We note that we were able to arrive at accurate and precise posterior estimates using the neural density estimator that simultaneously optimized an optimal design $\xi^*$, without any post-processing such as using SNL or Sampling Importance Resampling.

| Design Dimension | $\theta_0$ | $\theta_1$ |
| --- | --- | --- |
| D=1 | $1.29 \pm 2.98$ | $5.20 \pm 0.41$ |
| D=10 | $0.07 \pm 1.40$ | $4.87 \pm 0.16$ |
| D=100 | $1.35 \pm 0.52$ | $4.81 \pm 0.20$ |

*Table 1.* Posterior estimates mean and 68% interquartile range after observing $\xi^*$ values for each design dimension only using the amortized likelihood approximation provided by the neural density estimator used in the LF-PCE training. The held-out parameter values that were used to generate $y_o$ were $\boldsymbol{\theta}_{\text{true}} = [2, 5]$. More design dimensions approach the true held-out parameter with increasing precision.

## 5. Discussion

We demonstrated a novel information bounds, $I_{LF-PCE}$, to perform gradient-based BOED using black-box simulators present in many SBI applications and obtained lower bounds of the EIG on a toy model across a range of experimental design dimensions to showcase its scalability. Optimizing designs in SBI applications provides a valuable preconditioning step to typical sequential SBI methods such as SNL that are based on observed experimental designs. Sidestepping Bayesian optimization can also help to accelerate model testing and feedback from real-world data. Future work will examine the tradeoff between design diversity for improved entropy reduction and neural density estimator robustness, similar to the exploration and exploitation tradeoff present in Bayesian optimization.

# References

Cannon, P., Ward, D., and Schmon, S. M. Investigating the impact of model misspecification in neural simulation-based inference. doi: 10.48550/arxiv.2209.01845. URL https://arxiv.org/abs/2209.01845v1.

Cranmer, K., Brehmer, J., and Louppe, G. The frontier of simulation-based inference. *Proceedings of the National Academy of Sciences*, 117(48):30055–30062, November 2020. ISSN 0027-8424. doi: 10.1073/pnas.1912789117. URL http://arxiv.org/abs/1911.01429. arXiv: 1911.01429 Publisher: Proceedings of the National Academy of Sciences.

Durkan, C., Murray, I., and Papamakarios, G. On contrastive learning for likelihood-free inference, February 2020. URL http://arxiv.org/abs/2002.03712. arXiv: 2002.03712 Publication Title: arXiv.

Foster, A., Jankowiak, M., Bingham, E., Horsfall, P., Teh, Y. W., Rainforth, T., and Goodman, N. Variational Bayesian optimal experimental design. *arXiv*, March 2019a. ISSN 23318422. URL http://arxiv.org/abs/1903.05480. arXiv: 1903.05480 Publisher: arXiv.

Foster, A., Jankowiak, M., O'Meara, M., Teh, Y. W., and Rainforth, T. A unified stochastic gradient approach to designing Bayesian-optimal experiments. *arXiv*, November 2019b. ISSN 23318422. URL http://arxiv.org/abs/1911.00294. arXiv: 1911.00294 Publisher: arXiv.

Greenberg, D., Nonnenmacher, M., and Macke, J. Automatic posterior transformation for likelihood-free inference. In *International Conference on Machine Learning*, pp. 2404–2414. PMLR, 2019.

Hermans, J., Delaunoy, A., Rozet, F., Wehenkel, A., Begy, V., and Louppe, G. A crisis in simulation-based inference? beware, your posterior approximations can be unfaithful. *Transactions on Machine Learning Research*.

Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

Kleinegesse, S. and Gutmann, M. U. Efficient bayesian experimental design for implicit models. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 476–485. PMLR, 2019.

Kleinegesse, S. and Gutmann, M. U. Bayesian experimental design for implicit models by mutual information neural estimation, February 2020. URL http://arxiv.org/abs/2002.08129. arXiv: 2002.08129 Publication Title: arXiv.

Kleinegesse, S. and Gutmann, M. U. Gradient-based Bayesian Experimental Design for Implicit Models using Mutual Information Lower Bounds. May 2021. URL https://arxiv.org/abs/2105.04379v1. arXiv: 2105.04379.

Lindley, D. V. *Bayesian statistics, a review*, volume 2. SIAM, 1972.

Lueckmann, J., Bassetto, G., Karaletsos, T., and Macke, J. Likelihood-free inference with emulator networks. arxiv e-prints. *arXiv preprint arXiv:1805.09294*, 2018.

Lueckmann, J.-M., Boelts, J., Greenberg, D., Goncalves, P., and Macke, J. Benchmarking simulation-based inference. In *International Conference on Artificial Intelligence and Statistics*, pp. 343–351. PMLR, 2021.

Oord, A. v. d., Li, Y., and Vinyals, O. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.

Papamakarios, G. and Murray, I. Fast e-free inference of simulation models with Bayesian conditional density estimation. In *Advances in Neural Information Processing Systems*, pp. 1036–1044, May 2016. URL http://arxiv.org/abs/1605.06376. arXiv: 1605.06376 Issue: Nips ISSN: 10495258.

Papamakarios, G., Sterratt, D. C., and Murray, I. Sequential neural likelihood: Fast likelihood-free inference with autoregressive flows, May 2018. URL http://arxiv.org/abs/1805.07226. arXiv: 1805.07226 Publication Title: arXiv.

Papamakarios, G., Nalisnick, E., Rezende, D. J., Mohamed, S., and Lakshminarayanan, B. Normalizing Flows for Probabilistic Modeling and Inference. December 2019. URL http://arxiv.org/abs/1912.02762. arXiv: 1912.02762.

Poole, B., Ozair, S., Van Den Oord, A., Alemi, A., and Tucker, G. On variational bounds of mutual information. In *International Conference on Machine Learning*, pp. 5171–5180. PMLR, 2019.

Sisson, S., Fan, Y., and Beaumont, M. Overview of approximate bayesian computation. arxiv e-prints, art. *arXiv preprint arXiv:1802.09720*, 2018.

Talts, S., Betancourt, M., Simpson, D., Vehtari, A., and Gelman, A. Validating bayesian inference algorithms with

simulation-based calibration. doi: 10.48550/arxiv.1804. 06788. URL https://arxiv.org/abs/1804. 06788v2.

## A. Design Gradients of LF-PCE

For LF-PCE, we need unbiased gradient estimators of the information bounds. A normalizing flow can be seen as a reparameterized distribution, which allows for calculating the gradient with respect to designs $\nabla_\xi \boldsymbol{f}^{-1}(u; \theta, \xi)$. In practice, since we are evaluating the log probability of a data point, we would actually evaluate the inverse direction of a flow $\nabla_\xi \boldsymbol{f}(y; \theta, \xi)$ at the base distribution $p_u(\mathbf{u})$, which is usually a Gaussian distribution and evaluated by maximum likelihood.

More formally, following equation 4, the gradient with respect to $\xi$ is

$$\nabla_\xi \mathcal{L}(\xi) \approx -\frac{1}{N} \nabla_\xi \sum_n \log p_u(\boldsymbol{f}^{-1}(\mathbf{y}_n; \phi, \theta, \xi) + \log|\det \boldsymbol{J}(\boldsymbol{f}^{-1})(\mathbf{y}_n; \phi, \theta, \xi)|), \tag{10}$$

which is tractable as long as we can compute $\boldsymbol{f}^{-1}$, its Jacobian determinant, and evaluate the base density, $p_u(u)$, which is tractable for a base Gaussian distribution. Given this gradient, we can plug this into the gradient of LF-PCE to estimate the gradient of the information bound:

$$\frac{\partial I_{LF-PCE}}{\partial \xi} = \mathbb{E}_{p(\theta_0)p(y|\theta,\xi)q(\theta_{1:L}|y)} \left[ \frac{\partial g}{\partial \xi} + g \cdot \frac{\partial}{\partial \xi} \log p\phi(y|\theta_0, \xi) \right], \tag{11}$$

where

$$g(y, \theta_{0:L}, \phi, \xi) = \log \frac{p_\phi(y|\theta_0, \xi)}{\frac{1}{L+1}\sum_{\ell=0}^{L} p_\phi(y|\theta_\ell, \xi)}. \tag{12}$$

## B. Evaluation of Linear Model Designs and Posteriors

We evaluated the efficacy of the neural density estimator trained using the LF-PCE loss function to infer a held out true parameter value in Figure 2 by MCMC. We provide a quantitative evaluation of the posteriors in Table 1. The posteriors can be improved by computationally efficient methods such as Sampling Importance Resampling, or used in SBI algorithms that use sequential methods to refine the neural density estimator.



*Figure 2.* Comparison of the prior density the posterior achieved by the different design dimensional normalizing flows evaluated at an optimal design $p(\theta|y_o, \xi^*) \propto p_\phi(y_o|\theta, \xi^*)p(\theta)$. The red cross denotes the true model parameters.

## C. Evaluation of Posterior Predictive Distribution

As a reference, we plot the prior and posterior predictive plots for the 1-dimensional optimal design in Figure 3. An insight into the optimal experimental design problem is that the designs closer to where the prior distribution has more noise will lead to more clarification in a performed experiment, which is why the most optimal designs will be at the boundaries for the noisy linear model.

*Figure 3.* Prior predictive (blue) and posterior predictive (orange) distributions with the ground truth liner model (dotted red) for the single design case where D=1.