

---

# Implicit Deep Adaptive Design: Policy-Based Experimental Design without Likelihoods

---

Desi R. Ivanova<sup>†</sup> Adam Foster<sup>†</sup> Steven Kleinegesse<sup>‡</sup> Michael U. Gutmann<sup>‡</sup> Tom Rainforth<sup>†</sup>

<sup>†</sup>Department of Statistics, University of Oxford

<sup>‡</sup>School of Informatics, University of Edinburgh

`desi.ivanova@stats.ox.ac.uk`

## Abstract

We introduce implicit Deep Adaptive Design (iDAD), a new method for performing adaptive experiments in *real-time* with *implicit* models. iDAD amortizes the cost of Bayesian optimal experimental design (BOED) by learning a design *policy network* upfront, which can then be deployed quickly at the time of the experiment. The iDAD network can be trained on any model which simulates differentiable samples, unlike previous design policy work that requires a closed form likelihood and conditionally independent experiments. At deployment, iDAD allows design decisions to be made in milliseconds, in contrast to traditional BOED approaches that require heavy computation during the experiment itself. We illustrate the applicability of iDAD on a number of experiments, and show that it provides a fast and effective mechanism for performing adaptive design with implicit models.

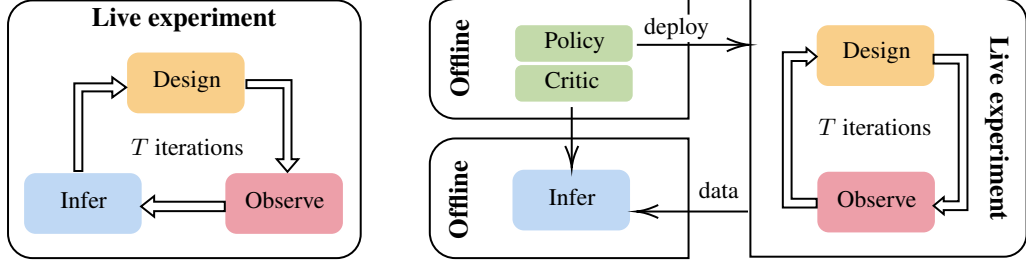
## 1 Introduction

Designing experiments to maximize the information gathered about an underlying process is a key challenge in science and engineering. Most such experiments are naturally *adaptive*—we can design later iterations on the basis of data already collected, refining our understanding of the process with each step [36, 45, 51]. For example, suppose that a chemical contaminant has accidentally been released and is rapidly spreading; we need to quickly discover its unknown source. To this end, we measure the contaminant concentration level at locations  $\xi_1, \dots, \xi_T$  (our experimental designs), obtaining observations  $y_1, \dots, y_T$ . Provided we can perform the necessary computations sufficiently quickly, we can design each  $\xi_t$  using data from steps  $1, \dots, t-1$  to narrow in on the source.

Bayesian optimal experimental design (BOED) [7, 32] is a principled model-based framework for choosing designs optimally; it has been successfully adopted in a diverse range of scientific fields [52, 58, 60]. In BOED, the unknown quantity of interest (e.g. contaminant location) is encapsulated by a parameter  $\theta$ , and our initial information about it by a prior  $p(\theta)$ . A simulator, or likelihood, model  $y|\theta, \xi$  describes the relationship between  $\theta$ , our controllable design  $\xi$ , and the experimental outcome  $y$ . To select designs *optimally*, the guiding principle is *information maximization*—we select the design that maximizes the expected (Shannon) information gained about  $\theta$  from the data  $y$ , or, equivalently, that maximizes the mutual information between  $\theta$  and  $y$ .

This naturally extends to adaptive settings by considering the *conditional* expected information gain given previously collected data. The traditional approach, depicted in Figure 1a, is to fit a posterior  $p(\theta|\xi_{1:t-1}, y_{1:t-1})$  after each iteration, and then select  $\xi_t$  in a myopic fashion using the one-step mutual information (see, e.g., [51] for a review). Unfortunately, this approach necessitates significant computation at each  $t$  and does not lend itself to selecting optimal designs quickly and adaptively.

Recently, Foster et al. [17] proposed an exciting alternative approach, called Deep Adaptive Design (DAD), that is based on learning design *policies*. DAD provides a way to avoid significant computation



(a) Traditional BOED: costly computations (design optimisation and parameter inference) are required at each iteration. (b) Policy-based BOED using iDAD: a design policy and critic are learnt before the live experiment. The policy enables quick and adaptive experiments, the critic assists likelihood-free inference.

Figure 1: Overview of adaptive BOED approaches applicable to implicit models.

at deployment-time by, prior to the experiment itself, learning a design policy network that takes past design-outcome pairs and near-instantaneously returns the design for the next stage of the experiment. The required training is done using simulated experimental histories, without the need to estimate any posterior or marginal distributions. DAD further only needs a single policy network to be trained for multiple experiments, further allowing for *amortization* of the adaptive design process. Unfortunately, DAD requires conditionally independent experiments and only works for the restricted class of models that have an explicit likelihood model we can simulate from, evaluate the density of, and calculate derivatives for, substantially reducing its applicability.

To address this shortfall, we instead consider a far more general class of models where we require only the ability to simulate  $y|\theta, \xi$  and compute the derivative  $\partial y/\partial \xi$ , e.g. via automatic differentiation [5]. Such models are ubiquitous in scientific modelling and include differentiable *implicit models* [19], for which the likelihood density  $p(y|\theta, \xi)$  is intractable. Examples include mixed effects models [15, 18], various models from chemistry and epidemiology [1], the Lotka Volterra model used in ecology [19], and models specified via stochastic differential equations (such as the SIR model [10]).

To perform rapid adaptive experimentation with this large class of models, we introduce *implicit Deep Adaptive Design* (iDAD), a method for learning adaptive design policy networks using only simulated outcomes (see Figure 1b). To achieve this, we introduce likelihood-free lower bounds on the total information gained from a sequence of experiments, which iDAD utilizes to learn a deep policy network. This policy network amortizes the cost of experimental design for implicit models and can be run in milliseconds at deployment-time. To train it, we show how the InfoNCE [57] and NWJ [37] bounds, popularized in representation learning, can be applied to the policy-based experimental design setting. The optimization of both of these bounds involves simultaneously learning an auxiliary *critic network*, bringing an important added benefit: it can be used to perform likelihood-free posterior inference of the parameters given the data acquired from the experiment.

We also relax DAD’s requirement for experiments to be conditionally independent, allowing its application in complex settings like time series data, and, through innovative architecture adaptations, also provide improvements in the conditionally independent setting as well. This further expands the model space for policy-based BOED, and leads to additional performance improvements.

Critically, iDAD forms the first method in the literature that can practically perform real-time adaptive BOED with implicit models: previous approaches are either not fast enough to run in real-time for non-trivial models, or require explicit likelihood models. We illustrate the applicability of iDAD on a range of experimental design problems, highlighting its benefits over existing baselines, even finding that it often outperforms costly non-amortized approaches. Code for iDAD is publicly available at <https://github.com/desi-ivanova/idad>.

## 2 Background

The BOED framework [32] begins by specifying a Bayesian model of the experimental process, consisting of a prior on the unknown parameters  $p(\theta)$ , a set of controllable designs  $\xi$ , and a data generating process that depends on them  $y|\theta, \xi$ ; as usual in BOED, we assume that  $p(\theta)$  does not depend on  $\xi$ . In this paper, we consider the situation where  $y|\theta, \xi$  is specified *implicitly*. This means

that it is defined by a deterministic transformation,  $f(\varepsilon; \theta, \xi)$ , of a base (or noise) random variable,  $\varepsilon$ , that is independent of the parameters and the design; e.g.,  $\varepsilon \sim \mathcal{N}(\varepsilon; 0, I)$ . The function  $f$  is itself often not known explicitly in closed form, but is implemented as a stochastic computer program (i.e. simulator) with input  $(\theta, \xi)$  and  $\varepsilon$  corresponding to the draws from the underlying random number generator (or equivalently the random seed). Regardless, the resulting induced likelihood density  $p(y|\theta, \xi)$  is still generally intractable, but sampling  $y|\theta, \xi$  is possible.

Having acquired a design-outcome pair  $(\xi, y)$ , we can quantify the amount of information we have gained about  $\theta$  by calculating the reduction in entropy from the prior to the posterior. We can further assess the quality of a design  $\xi$  before acquiring  $y$ , by computing the expected reduction in entropy with respect to the marginal distribution of the outcome,  $p(y|\xi) = \mathbb{E}_{p(\theta)}[p(y|\theta, \xi)]$ . The resulting quantity, called the *expected information gain* (EIG), is of central interest in BOED and is defined as

$$I(\xi) := \mathbb{E}_{p(\theta)p(y|\theta, \xi)} \left[ \log \frac{p(\theta|\xi, y)}{p(\theta)} \right] = \mathbb{E}_{p(\theta)p(y|\theta, \xi)} \left[ \log \frac{p(y|\theta, \xi)}{p(y|\xi)} \right]. \quad (1)$$

Note that  $I(\xi)$  is equivalent to the mutual information (MI) between the parameters  $\theta$  and data  $y$  when performing experiment  $\xi$ . The optimal  $\xi$  is then the one that maximises the EIG, i.e.  $\xi^* = \arg \max_{\xi} I(\xi)$ . Performing this optimization is a major computational challenge since the information objective is doubly intractable [46]. For implicit models, the cost becomes even greater as the likelihood is also not available in closed form, so estimating it, along with the marginal likelihood  $p(y|\xi)$ , is already itself a major computational problem [11, 20, 33, 54].

Jointly optimizing the design variables for all undertaken experiments at the same time using (1) is called *static* experimental design. In practice, however, we are often more interested in performing multiple experiments *adaptively* in a sequence  $\xi_1, \dots, \xi_T$ , so that the choice of each  $\xi_t$  can be guided by past experiments, namely the corresponding *history*  $h_{t-1} := \{(\xi_i, y_i)\}_{i=1:t-1}$ . The typical approach in such settings is to sequentially perform (approximate) posterior inference for  $\theta|h_{t-1}$ , followed by a one-step look ahead (myopic) BOED optimization that conditions on the observed history. In other words, to determine the designs  $\xi_1, \dots, \xi_T$ , we sequentially optimize the objectives

$$I_{h_{t-1}}(\xi_t) := \mathbb{E}_{p(\theta|h_{t-1})p(y_t|\theta, \xi_t, h_{t-1})} \left[ \log \frac{p(y_t|\theta, \xi_t, h_{t-1})}{p(y_t|\xi_t, h_{t-1})} \right], \quad t = 1, \dots, T. \quad (2)$$

However, such approaches incur significant computational cost during the experiment itself, particularly for implicit models [16, 21, 30]. This has critical consequences: in most cases they cannot be run in real-time, undermining one's ability to use them in practice.

## 2.1 Policy-based adaptive design with likelihoods

For tractable likelihood models, Foster et al. [17] proposed a new framework, called Deep Adaptive Design (DAD), for adaptive experimental design that avoids expensive computations during the experiment. To achieve this, they introduce a parameterized deterministic design function, or policy,  $\pi_\phi$  that takes the history  $h_{t-1}$  as input and returns the design  $\xi_t = \pi_\phi(h_{t-1})$  to be used for the next experiment as output. This set-up allows them to consider the objective

$$\mathcal{I}_T(\pi_\phi) = \mathbb{E}_{p(\theta)p(h_T|\theta, \pi_\phi)} \left[ \sum_{t=1}^T I_{h_{t-1}}(\xi_t) \right], \quad \xi_t = \pi_\phi(h_{t-1}), \quad (3)$$

which crucially depends on the policy  $\pi$  rather than the individual design  $\xi_t$ . Learning a policy up-front, rather than designs, is what allows adaptive experiments to be performed in real-time.

Under the assumption that  $y_t$  is independent of  $h_{t-1}$  conditional on the parameters  $\theta$  and the design  $\xi_t$ , i.e.  $p(y_t|\theta, \xi_t, h_{t-1}) = p(y_t|\theta, \xi_t)$ , Foster et al. [17] showed that the objective can be simplified to

$$\mathcal{I}_T(\pi_\phi) = \mathbb{E}_{p(\theta)p(h_T|\theta, \pi_\phi)} \left[ \log \frac{p(h_T|\theta, \pi_\phi)}{p(h_T|\pi_\phi)} \right], \quad p(h_T|\theta, \pi_\phi) = \prod_{t=1}^T p(y_t|\theta, \xi_t). \quad (4)$$

To deal with the marginal  $p(h_T|\pi_\phi)$  in the denominator, they then derived several optimizable lower bounds on  $\mathcal{I}_T(\pi_\phi)$ , such as the sequential Prior Contrastive Estimation (sPCE) bound

$$\mathcal{L}_T^{\text{sPCE}}(\pi_\phi, L) = \mathbb{E}_{p(\theta_0)p(h_T|\theta, \pi_\phi)p(\theta_{1:L})} \left[ \log \frac{p(h_T|\theta_0, \pi_\phi)}{\frac{1}{L+1} \sum_{\ell=0}^L p(h_T|\theta_\ell, \pi_\phi)} \right] \leq \mathcal{I}_T(\pi_\phi) \quad \forall L \geq 1. \quad (5)$$

The parameters of the policy  $\phi$ , which takes the form of a deep neural network, are now learned prior to the experiment(s) using stochastic gradient ascent on this bound with simulated experimental histories. Design decisions can then be made using a single forward pass of  $\pi_\phi$  during deployment. Unfortunately, training the DAD network by optimizing (5) requires the likelihood density  $p(h_T|\theta, \pi)$  to be analytically available—an assumption that is too restrictive in many practical situations. The architecture for DAD also assumes conditionally independent designs, which is unsuitable in some settings like time-series data. Our method lifts both of these restrictions.

### 3 Implicit Deep Adaptive Design

We have seen that the traditional step-by-step approach to adaptive design for implicit models [16, 21, 30] is too costly to deploy for most applications, whilst the only existing policy-based approach, DAD [17], makes restrictive assumptions that prevent it being applied to implicit models. We aim to relax the restrictive assumptions of the latter, making policy-based BOED applicable to all models where we can sample from  $y|\theta, \xi$  and compute the derivative  $\partial y/\partial \xi$ , a strict superset of the class of models that can be handled by DAD. This requires new training objectives for the policy network that do not involve an explicit likelihood and are not based on conditionally independent experiments, along with new architectures that work for non-exchangeable models like time series.

#### 3.1 Information lower bounds for policy-based experimental design without likelihoods

To establish a suitable likelihood-free training objective for the implicit setting, our high-level idea is to leverage recent advances in variational MI [see 42, for an overview], which have shown promise for *static* BOED [16, 28, 29]. While using these bounds in the traditional framework of (2) would not permit real-time experiments, one could consider a naive application of them to the policy objective of (3) by replacing each  $I_{h_{t-1}}$  with a suitable variational lower bound that uses a ‘critic’  $U_t : \mathcal{H}^{t-1} \times \Theta \rightarrow \mathbb{R}$  to avoid explicit likelihood evaluations, where  $\mathcal{H}^{t-1}$  and  $\Theta$  are the spaces of histories and parameters respectively. An effective critic successfully encapsulates the true likelihood, tightening the bound. Although its form depends on the choice of bound, all critics are parametrized and trained in the same way, namely by a neural network  $U_{\phi_t}$  which is optimized to tighten the bound. Unfortunately, replacing each  $I_{h_{t-1}}$  involves learning  $T$  such critic networks and requires samples from all posteriors  $p(\theta|h_{t-1})$ , which will typically be impractically costly.

To avoid this issue, we show that we can obtain a unified information objective similar to (4), *even without conditionally independent experiments*. The following proposition therefore marks the first key milestone in eliminating the restrictive assumptions of [17], by establishing a unified objective without intermediate posteriors that is valid even when the model itself changes between time steps.

**Proposition 1** (Generalized total expected information gain). *Consider the data generating distribution  $p(h_T|\theta, \pi) = \prod_{t=1:T} p(y_t|\theta, \xi_t, h_{t-1})$ , where  $\xi_t = \pi(h_{t-1})$  are the designs generated by the policy and, unlike in (4),  $y_t$  is allowed to depend on the history  $h_{t-1}$ . Then we can write (3) as*

$$\mathcal{I}_T(\pi) = \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} [\log p(h_T|\theta, \pi)] - \mathbb{E}_{p(h_T|\pi)} [\log p(h_T|\pi)]. \quad (6)$$

Proofs are presented in Appendix A. The advantage of (6) is that we can draw samples from  $p(\theta)p(h_T|\theta, \pi)$  simply by sampling our model and taking forward passes through the design network. However, neither of the *densities*  $p(h_T|\theta, \pi)$  and  $p(h_T|\pi)$  are tractable for implicit models.

To side-step this intractability, we observe that  $\mathcal{I}_T(\pi)$  takes an analogous form to a MI between  $\theta$  and  $h_T$ . For measure-theoretic reasons, namely because the  $\xi_{1:T}$  are deterministic given  $y_{1:T}$  (see Appendix A for a full discussion), it is not the true MI. However, the following two propositions show that we can treat  $\mathcal{I}_T(\pi)$  *as if it were this MI*. Specifically, we show that the InfoNCE [57] and NWJ [37] bounds on the MI can be adapted to establish tractable lower bounds on our unified objective  $\mathcal{I}_T(\pi)$ . These two bounds both utilize a *single* auxiliary critic network  $U_\psi$  that is trained simultaneously with the design network.

**Proposition 2** (NWJ bound for implicit policy-based BOED). *For a design policy  $\pi$  and a critic function  $U : \mathcal{H}^T \times \Theta \rightarrow \mathbb{R}$ , let*

$$\mathcal{L}_T^{NWJ}(\pi, U) := \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} [U(h_T, \theta)] - e^{-1} \mathbb{E}_{p(\theta)p(h_T|\pi)} [\exp(U(h_T, \theta))], \quad (7)$$

*then  $\mathcal{I}_T(\pi) \geq \mathcal{L}_T^{NWJ}(\pi, U)$  holds for any  $U$ . Further, the inequality is tight for the optimal critic  $U_{NWJ}^*(h_T, \theta) = \log p(h_T|\theta, \pi) - \log p(h_T|\pi) + 1$ .*

---

**Algorithm 1:** Implicit Deep Adaptive Design with (iDAD)

---

**Input:** Differentiable simulator  $f$ , sampler for prior  $p(\theta)$ , number of experimental steps  $T$

**Output:** Design network  $\pi_\phi$ , critic network  $U_\psi$

**while** Computational training budget not exceeded **do**

    Sample  $\theta \sim p(\theta)$  and set  $h_0 = \emptyset$

**for**  $t = 1, \dots, T$  **do**

        Compute  $\xi_t = \pi_\phi(h_{t-1})$

        Sample  $\varepsilon_t \sim p(\varepsilon)$  and compute  $y_t = f(\varepsilon_t; \xi_t, \theta, h_{t-1})$

        Set  $h_t = \{(\xi_1, y_1), \dots, (\xi_t, y_t)\}$

**end**

    Estimate  $\nabla_{\phi, \psi} \mathcal{L}_T(\pi_\phi, U_\psi)$  as per (10) where  $\mathcal{L}_T$  is  $\mathcal{L}_T^{\text{NWJ}}$  (7) or  $\mathcal{L}_T^{\text{NCE}}$  (8)

    Update the parameters  $(\phi, \psi)$  using stochastic gradient ascent scheme

**end**

For deployment, use the deterministic trained design network  $\pi_\phi$  to obtain a designs  $\xi_t$  directly.

---

**Proposition 3** (InfoNCE bound for implicit policy-based BOED). *Let  $\theta_{1:L} \sim p(\theta_{1:L}) = \prod_i p(\theta_i)$  be a set of contrastive samples where  $L \geq 1$ . For design policy  $\pi$  and critic function  $U: \mathcal{H}^T \times \Theta \rightarrow \mathbb{R}$ , let*

$$\mathcal{L}_T^{\text{NCE}}(\pi, U; L) := \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)} \mathbb{E}_{p(\theta_{1:L})} \left[ \log \frac{\exp(U(h_T, \theta_0))}{\frac{1}{L+1} \sum_{i=0}^L \exp(U(h_T, \theta_i))} \right], \quad (8)$$

*then  $\mathcal{I}_T(\pi) \geq \mathcal{L}_T^{\text{NCE}}(\pi, U; L)$  for any  $U$  and  $L \geq 1$ . Further, the optimal critic,  $U_{\text{NCE}}^*(h_T, \theta) = \log p(h_T|\theta, \pi) + c(h_T)$  where  $c(h_T)$  is any arbitrary function depending only on the history, recovers the sPCE bound in (5); the inequality is tight in the limit as  $L \rightarrow \infty$  for this optimal critic.*

We propose these two alternative bounds due to their complementary properties: the NWJ bound can have large variance, but tends to be less biased. That is, the NWJ bound tends to be tighter for good critics, but is itself more difficult to reliably estimate and thus optimize. While the NWJ critic must learn to self-normalize, the InfoNCE bound avoids this issue but typically will not be tight for finite  $L$  even with an optimal critic (note  $\mathcal{L}_T^{\text{NCE}} \leq \log(L+1)$  [42]). Consequently, only the NWJ objective recovers the true optimal policy if our critic has infinite capacity and our optimization scheme is perfect, i.e.  $\arg \max_\pi \max_U \mathcal{L}_T^{\text{NWJ}}(\pi, U) = \pi^* \neq \arg \max_\pi \max_U \mathcal{L}_T^{\text{NCE}}(\pi, U; L)$  in general, but it can be more difficult to work with in practice. We present a third bound that provides a potential solution to this, and further discuss the relative merits of the two bounds, in Appendix A.

We note that for both bounds the optimal critic does not depend on the learned policy. The final trained critic can be used to approximate the density ratio  $p(h_T|\theta, \pi)/p(h_T|\pi) = p(\theta|h_T)/p(\theta)$ , either directly in the case of the NWJ critic, or via self-normalization for the InfoNCE bound. We can use this to help approximate the posterior over  $\theta$  given the collected real data from the experiment. This means we can perform likelihood-free inference after training the critic, which extends previous results [28, 29] from the static to the adaptive policy-based setting.

### 3.2 Parameterization and gradient estimation

In practice, we represent the policy  $\pi$  and the critic  $U$  as neural networks,  $\pi_\phi$  and  $U_\psi$  respectively, such that the lower bounds become a function  $\mathcal{L}(\pi_\phi, U_\psi)$  of their parameters. By simultaneously optimizing  $\mathcal{L}(\pi_\phi, U_\psi)$  with respect to both  $\phi$  and  $\psi$ , we both learn a tight bound that accurately represents the true MI and a design policy network that produces high-quality designs under this metric.

We optimize these bounds using stochastic gradient methods [26, 49]. For this, we must account for the fact that the parameter  $\phi$  affects the probability distributions with respect to which expectations are taken. We deal with this problem by utilizing the reparametrization trick [35, 48], for which we assume that design space  $\Xi$  and observation space  $\mathcal{Y}$  are continuous. To this end, we first formalize the notion of a differentiable implicit model in the adaptive design setting as

$$y_t = f(\varepsilon_t; \xi_t(h_{t-1}), \theta, h_{t-1}), \quad \text{where } \theta \sim p(\theta), \quad \varepsilon_t \sim p(\varepsilon) \quad \forall t \in \{1, \dots, T\} \quad (9)$$

and we assume that we can compute the derivatives  $\partial f / \partial \xi$  and  $\partial f / \partial h$ . Interestingly, it is possible to use an implicit prior without access to the density  $p(\theta)$ , and we do not need access to  $\partial f / \partial \theta$ .



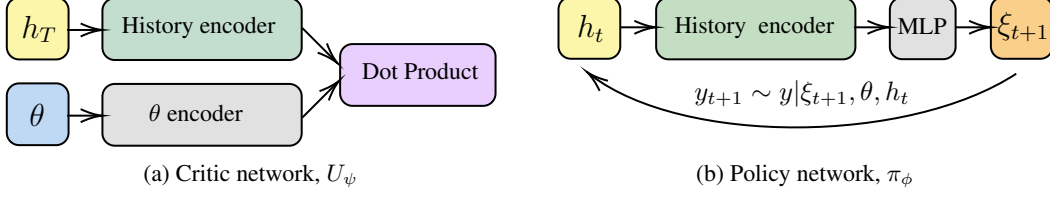


Figure 2: Overview of network architectures used in iDAD.

Under these conditions, we can express the bounds in terms of expectations that do not depend on  $\phi$  or  $\psi$ , and hence move the gradient operator inside. For  $\mathcal{L}_T^{\text{NCE}}(\pi_\phi, U_\psi; L)$ , for example, we have

$$\nabla_{\phi, \psi} \mathcal{L}_T^{\text{NCE}} = \mathbb{E}_{p(\theta_{0:L})p(\varepsilon_{1:T})} \left[ \nabla_{\phi, \psi} \log \frac{\exp(U_\psi(h_T(\varepsilon_{1:T}, \pi_\phi), \theta_0))}{\frac{1}{L+1} \sum_{i=0}^L \exp(U_\psi(h_T(\varepsilon_{1:T}, \pi_\phi), \theta_i))} \right]. \quad (10)$$

While each element of the history  $h_T$  depends on  $\phi$  in a possibly nested manner, we do not need to explicitly keep track of these dependencies thanks to automatic differentiation [5, 41].

Like DAD, our new method—which we call *implicit Deep Adaptive Design* (iDAD)—is trained with simulated histories  $h_T = \{(\xi_i, y_i)\}_{i=1:T}$  prior to the actual experiment, allowing design decision to be made using a single forward pass during deployment. Unlike DAD, however, it does not require knowledge of the likelihood function, nor the assumption of conditionally independent designs, which significantly broadens its applicability. A summary of the iDAD approach is given in Algorithm 1.

### 3.3 Network architectures

The iDAD approach involves the simultaneous training of the *policy*  $\pi_\phi$  and *critic*  $U_\psi$  networks. It is essential to choose the neural architectures of these two components carefully to learn effective policies: poor choices of critic architecture will lead to loose, unrepresentative, bounds, while poor choices of policy architecture will directly lead to ineffective policies. Good choices of architecture need to balance flexibility with ease of training, and will typically require the incorporation of problem-specific inductive biases. A high-level summary of our architectures is shown in Figure 2.

The critic network,  $U_\psi$ , takes a *complete* history  $h_T$  and the parameter  $\theta$  as input, and outputs a scalar. Our suggested architecture first encodes the two inputs separately to representations of the same dimension, using a *history encoder*,  $E_{\psi_h}$ , and a *parameter encoder*,  $E_{\psi_\theta}$ , respectively. The output of the critic is then simply taken as their dot product  $U_\psi(h_T, \theta) := E_{\psi_h}(h_T)^\top E_{\psi_\theta}(\theta)$ ; after training, the two encodings correspond to approximate sufficient statistics [9]. This setup corresponds to a separable critic architecture, as is commonly used in the representation learning literature [3, 8, 57]. While we use a simple MLP for  $E_{\psi_\theta}$ , the setup for  $E_{\psi_h}$  varies with the context as we discuss below.

The policy network,  $\pi_\phi$ , takes the available history,  $h_t$ , as input, and outputs a design. Our suggested architecture makes use of a history encoder,  $E_{\phi_h}$ , of the same form as  $E_{\psi_h}$ , except that it must now take in varying length inputs; its output remains a fixed dimensional embedding. We then pass this embedding through an MLP to produce the design  $\xi_{t+1}$ . At the next iteration of the experiment, the same policy network is then called again with the updated history  $h_{t+1} = h_t \cup \{(\xi_{t+1}, y_{t+1})\}$ .

We use the same architecture for both history encoders,  $E_{\psi_h}$  and  $E_{\phi_h}$ , but do not share network parameters between them. This architecture first individually embeds each design–outcome pair  $(\xi_t, y_t)$  to a corresponding representation,  $r_t$ , using a simple MLP that is shared across all time steps. The produced history encoding is then an aggregation of these representations, with how this is done depending on whether the experiments are *conditionally independent*, i.e.  $y_t \perp\!\!\!\perp h_{t-1} | \theta, \xi_t$ , or not.

Foster et al. [17] proved that if experiments are conditionally independent, then the optimal policy is invariant to the order of the history. We prove that the same is true for the critic in Proposition 5 in Appendix C. In our setup, we can exploit this result by using a *permutation invariant* aggregation strategy for  $\{r_1, \dots, r_t\}$  when conditional independence holds. The simplest approach to do this would be to use sum-pooling [62], as was done in DAD. However, to improve on this, we instead propose using a more advanced permutation invariant architecture based on self-attention [13, 25, 39, 47, 59], namely that of Parmar et al. [40]; we find this provides notable empirical gains. When conditional independence does not hold, this approach is no longer appropriate and we instead use an LSTM [22] for the aggregation. See Appendix C for further details.

## 4 Related work

Adaptive policy-based BOED has only recently been introduced [17] and has not yet been extended to implicit models—the gap that this work addresses. Previous approaches to adaptive experiments usually follow the two-step greedy procedure described in Section 2. Methods for MI/EIG estimation without likelihoods include the use of variational bounds [15, 16, 28] and ratio estimation [27, 30]; approximate Bayesian computation together with kernel density estimation [43]; and approximating the intractable likelihood first, for example via polynomial chaos expansion [24], followed by applying likelihood-based estimators, such as nested Monte Carlo [46]. The maximization step in more traditional methods tends to rely on gradient-free optimization, including grid-search, evolutionary algorithms [44], Bayesian optimization [15, 30], or Gaussian process surrogates [38]. More recently, gradient-based approaches have been introduced [15, 28], some of which allow the estimation and optimization simultaneously in a single stochastic-gradient scheme [16, 23, 29]. From a posterior estimation perspective, likelihood-free inference can be performed via approximate Bayesian computation [33, 54], ratio estimation [56], conventional MCMC for methods that make tractable approximation to the likelihood [23, 24], or as a byproduct of MI estimation [16, 27, 29, 30].

## 5 Experiments

We evaluate the performance of **iDAD** on a number of real-world experimental design problems and a range of baselines. A summary of all the methods that we consider is given in Table 1. Since we aim to perform adaptive experiments in *real-time*, we focus mostly on baselines that do not require significant computational time during the experiment. These include heuristic approaches that require no training, namely **equal** interval designs (when possible) and **random** designs, as well as static BOED approaches, where we, non-adaptively, choose all the designs prior to the experiment by optimising the mutual information objective of Equation (1) with  $\xi = \{\xi_1, \dots, \xi_T\}$  and  $y = \{y_1, \dots, y_T\}$ . The static BOED approaches we consider are the **MINEBED** method of [28] and the likelihood-free ACE approach of [16], where we use the prior as a proposal distribution, referring to this baseline as **SG-BOED**. We also implement the expensive traditional non-amortized myopic strategy described in Section 2, for which we use the **variational** approach of [16], with the Barber-Agakov bound [4, 15], at each experiment step (see Appendix D.3 for details). Finally, where possible, we compare our method with DAD [17], in order to assess the performance gap that would arise if we had an analytic likelihood. This comparison is done primarily for evaluation purposes—because it has access to the likelihood density, DAD serves as an upper bound on the performance iDAD can achieve; one should use explicit likelihood methods whenever possible.

The main performance metric that we focus on is the total EIG,  $\mathcal{I}_T(\pi)$ , as given in (6). In cases where the likelihood is available, we estimate the  $\mathcal{I}_T(\pi)$  using the sPCE lower bound in (5) and its sister upper bound, the sequential Nested Monte Carlo bound [sNMC; 17]. To ensure that the bounds are tight, we evaluate them with a large number of contrastive samples, i.e.  $L \geq 10^5$ . Where the likelihood is truly intractable, we assess the iDAD strategy in a more qualitative manner by looking at the optimal designs and approximate posteriors. For the adaptive experiments, we further consider the deployment time (i.e. the time required to propose a design), which is a critical metric for our aims. All deployment times exclude the time needed to determine the first experiment as it can be computed up-front, during the training phase. Timings for training the policy itself are given in Appendix D.

The main performance metric that we focus on is the total EIG,  $\mathcal{I}_T(\pi)$ , as given in (6). In cases where the likelihood is available, we estimate the  $\mathcal{I}_T(\pi)$  using the sPCE lower bound in (5) and its sister upper bound, the sequential Nested Monte Carlo bound [sNMC; 17]. To ensure that the bounds are tight, we evaluate them with a large number of contrastive samples, i.e.  $L \geq 10^5$ . Where the likelihood is truly intractable, we assess the iDAD strategy in a more qualitative manner by looking at the optimal designs and approximate posteriors. For the adaptive experiments, we further consider the deployment time (i.e. the time required to propose a design), which is a critical metric for our aims. All deployment times exclude the time needed to determine the first experiment as it can be computed up-front, during the training phase. Timings for training the policy itself are given in Appendix D.

### 5.1 Location Finding

We first demonstrate our approach on the location finding experiment from [17]. Inspired by the acoustic energy attenuation model [53], this experiment involves finding the locations of multiple hidden sources, each emitting a signal with intensity that decreases according to the inverse-square law. The *total intensity*—a superposition of these signals—can be measured noisily at any location. The design problem is choosing where to measure the total signal in order to uncover the sources.

Table 1: Key properties of considered methods.

	Adaptive	Real-time	Implicit
Random	✗	N/A	✓
Equal interval	✗	N/A	✓
MINEBED	✗	N/A	✓
SG-BOED	✗	N/A	✓
Variational	✓	✗	✓
DAD	✓	✓	✗
<b>iDAD</b>	✓	✓	✓

Table 2: Lower bounds on the total information,  $\mathcal{I}_{10}(\pi)$ , for the location finding experiment in Section 5.1. The bounds were estimated using  $L = 5 \times 10^5$  contrastive samples. Errors indicate  $\pm 1$  standard errors estimated over 4096 histories (128 for variational). Corresponding upper bounds are given in Table 6 in Appendix D.

Method \ $\theta$ dim.	4D	6D	10D	20D
Random	$4.791 \pm 0.040$	$3.468 \pm 0.014$	$1.889 \pm 0.011$	$0.552 \pm 0.006$
MINEBED	$5.518 \pm 0.028$	$4.221 \pm 0.028$	$2.458 \pm 0.029$	$0.801 \pm 0.019$
SG-BOED	$5.547 \pm 0.028$	$4.215 \pm 0.030$	$2.454 \pm 0.029$	$0.803 \pm 0.019$
Variational	$4.639 \pm 0.144$	$3.625 \pm 0.165$	$2.181 \pm 0.151$	$0.669 \pm 0.097$
<b>iDAD (NWJ)</b>	<b><math>7.694 \pm 0.045</math></b>	$5.765 \pm 0.036$	<b><math>3.252 \pm 0.039</math></b>	<b><math>0.877 \pm 0.022</math></b>
<b>iDAD (InfoNCE)</b>	<b><math>7.750 \pm 0.039</math></b>	<b><math>5.986 \pm 0.037</math></b>	<b><math>3.251 \pm 0.039</math></b>	<b><math>0.871 \pm 0.020</math></b>
DAD	$7.967 \pm 0.034$	$6.300 \pm 0.030$	$3.337 \pm 0.039$	$0.937 \pm 0.022$

Table 3: Lower and upper bounds on MI  $\mathcal{I}_{10}(\pi)$  for different network architectures on location finding experiment using the InfoNCE bound. All estimates obtained as in Table 2.

Design	Critic	Lower bound	Upper bound
<b>Attention</b>	<b>Attention</b>	<b><math>7.750 \pm 0.039</math></b>	<b><math>7.863 \pm 0.043</math></b>
Attention	Pooling	$7.567 \pm 0.037$	$7.632 \pm 0.039$
Pooling	Attention	$7.398 \pm 0.040$	$7.470 \pm 0.042$
Pooling	Pooling	$7.135 \pm 0.034$	$7.192 \pm 0.041$

In Table 2 we can see that iDAD substantially outperforms all baselines including, perhaps surprisingly, the traditional (non-amortized) adaptive variational approach, despite its large computational cost shown Table 4. The poor performance of the variational approach is likely driven by the inability of the mean-field variational family to capture the highly non-Gaussian true posterior, highlighting the detrimental effect that wrong posteriors can have on determining optimal designs when using the traditional sequential BOED approach.

Table 2 further shows that the performance gap to the likelihood-based DAD method is small, even as the dimension of the design and parameter space grows. Though the information gained by all methods decreases with the dimensionality, this is to be expected: in higher dimensions it is inherently more difficult to infer the relative direction of the sources from observing their intensity. Overall, this experiment demonstrates that iDAD is able to learn near-optimal amortized design policies without likelihoods, while being run in milliseconds at deployment.

**Ablation: attention to history.** We next assess the benefit of utilizing our proposed more sophisticated permutation invariant architectures, compared to the simple pooling of [62] used in [17]. Our approach incorporates attention layers into both networks that we train. This leads us to four possible combinations of network architectures. Table 3 compares the efficacy of the resulting design policies and strongly suggests that incorporating attention mechanisms in either and/or both networks improves performance, with inclusion in the design network particularly important.

We preform further ablation studies to investigate and demonstrate important properties of our method, such as its scalability with the number of experiments  $T$ , stability between different training runs and performance to errors in the design network (introduced by not training the network to convergence). Results and discussion are provided in Appendix D.4.4.

## 5.2 Pharmacokinetic model

Our next experiment is taken from the pharmacokinetics literature and has been studied in other recent works on BOED for implicit models [28, 63]. Specifically, we consider the compartmental model

Table 4: Deployment time of adaptive methods in 2D, measured on a CPU. Errors were calculated on the basis of 10 runs.

Method	Deployment time (sec.)
Variational	$2256.0 \pm 1\%$
<b>iDAD (NWJ)</b>	$0.0167 \pm 2\%$
<b>iDAD (InfoNCE)</b>	$0.0168 \pm 2\%$
DAD	$0.0070 \pm 6\%$



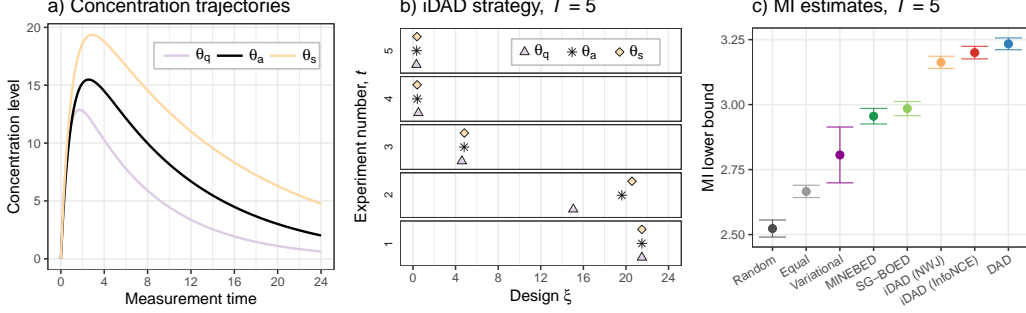


Figure 3: Plots for pharmacokinetics experiment. a) Visualisation of model showing concentration level as a function of measurement time for 3 values of  $\theta$ , resulting in a quick ( $\theta_q$ ), average ( $\theta_a$ ), or slow ( $\theta_s$ ) trajectory. b) Designs selected by an iDAD policy trained with InfoNCE. c) MI lower bounds achieved by iDAD and baselines. All estimates obtained as in Table 2.

of [50], for which the distribution of an administered drug through the body is governed by three parameters: absorption rate  $k_\alpha$ , elimination rate  $k_e$ , and volume  $V$ , which form the parameters of interest, i.e.  $\theta = (k_\alpha, k_e, V)$ . Given  $T = 5$  patients, the design problem is to adaptively choose blood sampling times,  $0 \leq \xi_t \leq 24$  hours, for each, measured from the the point the drug was administered (with patient 2 not being administered until after sampling patient 1 etc). Plausible concentration trajectories are shown in Figure 3a). Full details and further results are given in Appendix D.5.

We first qualitatively consider the design policy of iDAD (trained with the InfoNCE objective) in Figure 3b). As we have not yet observed any data, the optimal design for the first patient (bottom row) is the same for all  $\theta$ . For the second patient, only guided by  $\xi_1$  and the outcome  $y_1$ , iDAD is already able to distinguish between quickly and slowly decaying concentration trajectories: it proposes a significantly earlier measurement time for the quickly decaying trajectory (purple triangle,  $\theta_q$ ) and later time for the slowly decaying one (yellow diamond,  $\theta_s$ ). For the third patient, iDAD always targets the peak of the drug concentration distribution which is quite similar for all  $\theta$ . Measurements for the last two patients are made soon after the drug has been administered ( $\sim 15 - 30$  min), when concentration levels increase rapidly, to capture information about how quickly the drug is absorbed.

To provide more quantitative assessment and compare to our baselines, we again consider the final EIG values as shown in Figure 3c). This reveals that the iDAD strategies perform best among the methods that are applicable to implicit models, confirming that the learnt policies propose superior designs. The performance gap to DAD, which relies on explicit likelihoods, is not statistically significant (at the 5% level) for iDAD trained with InfoNCE, while significant, but still small, for NWJ.

Finally, we consider the convergence of the iDAD networks under the different training objectives and compare to DAD for reference. As shown in Figure 4, although all three converge to approximately the same value, they do so at rather different speeds: while DAD requires about 5000 gradient updates, implicit methods need longer training and tend to exhibit higher variance, particularly NWJ.

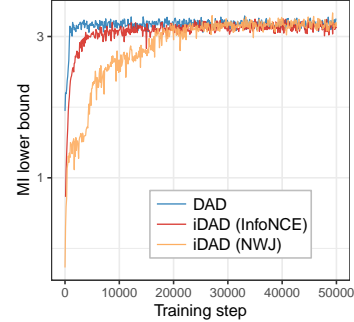


Figure 4: Convergence of MI lower bounds.

### 5.3 SIR Model

In this experiment, we demonstrate our approach on an implicit model from epidemiology. Namely, we consider a formulation of the stochastic SIR model [10] that is based on stochastic differential equations (SDEs), as done by [29]. Here, individuals in a fixed population belong to one of three categories: susceptible, infected or recovered. Susceptible people can become infected and then recover, with the dynamics of these two events being governed by two model parameters—the infection rate  $\beta$  and the recovery rate  $\gamma$ . Our aim is to determine the optimal times  $\tau$  at which to measure the number of infected people,  $I(\tau)$ , in order to estimate the two parameters. This implicit model is challenging because data simulation is expensive, since we need to solve many SDEs, and experimental designs have a time-dependency. See Appendix D.6 for full details.

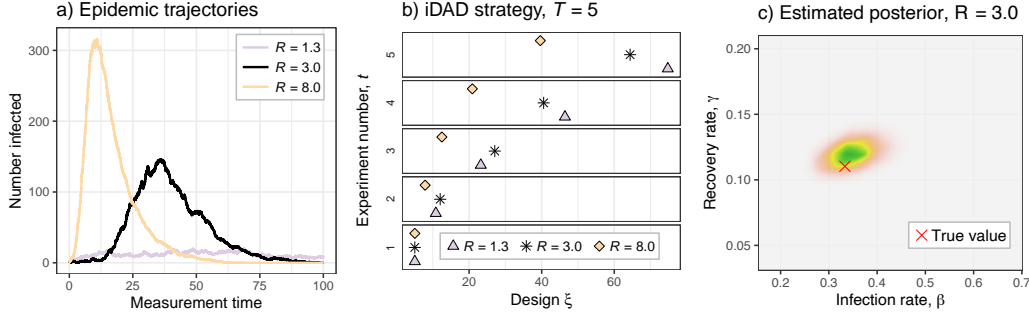


Figure 5: a) Epidemic trajectories for 3 realization of  $(\beta, \gamma)$  with different reproduction numbers  $R = \beta/\gamma$ . b) Designs selected by an iDAD policy trained with NWJ. c) Example posterior estimates from the critic network given data generated with the ground-truth parameters shown by the red cross.

We train a iDAD networks to perform  $T = 5$  experiments and compare against random, equal interval, and static design baselines; DAD cannot be run because the problem corresponds to a true implicit model. Table 5 shows lower bound estimates on the MI and demonstrates that iDAD outperforms all compared methods. Note that a degree of caution is required when analysing the results, as they are influenced by unavoidable biases in the estimation process. Namely, a critic is still required to estimate the MI lower bound, and there may be variations in the effectiveness of these critics, with less effective ones corresponding to looser bounds and therefore underestimating the true MI. Nonetheless, for other models where such checks are possible, we have found the bounds to be relatively tight, while, even if this turns out not to be the case here, the fact that the critics for the static approaches are easier to train should mean our relative evaluations for iDAD (and Random) are still conservative compared to the other baselines.

Table 5: MI lower bounds ( $\pm 1$  s.e.).

Method	Lower bound
Random	$1.915 \pm 0.032$
Equal interval	$2.669 \pm 0.023$
MINEBED	$3.400 \pm 0.001$
SG-BOED	$3.752 \pm 0.020$
<b>iDAD (NWJ)</b>	<b><math>3.869 \pm 0.001</math></b>
<b>iDAD (InfoNCE)</b>	<b><math>3.915 \pm 0.020</math></b>

Figure 5 further demonstrates important qualitative results for this model. Figure 5a) shows different epidemic trajectories, i.e. the number of infected  $I(\tau)$  people as a function of measurement time  $\tau$ , whilst 5b) plots their corresponding designs obtained from the learned iDAD policy. Importantly, diseases with a significantly different profile, e.g. a slow ( $R = 1.3$ ) or a fast ( $R = 8.0$ ) spread result in different sets of optimal designs, highlighting the adaptivity of iDAD. Finally, Figure 5c) shows an example posterior distribution estimate from the learnt iDAD critic network, which we see is consistent with the ground truth parameters.

## 6 Discussion

**Limitations.** The benefit that iDAD can be used in live experiments comes at the cost of substantial training that can be computationally expensive. However, this is mitigated by its amortization of the adaptive design process, such that only one network needs training, even if we have multiple experiment instances. The cost–performance trade-off can also be directly controlled by judicious choices of architecture and the amount of training performed. Another natural limitation is that the use of gradients naturally restricts the approach to continuous design settings, something which future work might look to address.

**Conclusions.** In this paper we introduced iDAD—the first policy-based adaptive BOED method that can be applied to implicit models. By training a design network without likelihoods upfront, iDAD is thus the first method that allows real-time adaptive experiments for simulator-based models. In our experiments, iDAD performed significantly better than all likelihood-free baselines. Further, by using models where the likelihood is available as a test bed, we found that it was able to almost match the analogous likelihood-based adaptive approach, which acts as an upper bound on what might be achieved without access to the likelihood itself. In conclusion, we believe iDAD marks a step change in Bayesian experimental design for *implicit* models, allowing designs to be proposed quickly, adaptively, and non-myopically during the live experiment.

## Acknowledgments and Disclosure of Funding

DRI is supported by EPSRC through the Modern Statistics and Statistical Machine Learning (StatML) CDT programme, grant no. EP/S023151/1. AF gratefully acknowledges funding from EPSRC grant no. EP/N509711/1. SK was supported in part by the EPSRC Centre for Doctoral Training in Data Science, funded by the UK Engineering and Physical Sciences Research Council (grant EP/L016427/1) and the University of Edinburgh.

## References

- [1] Edward J. Allen, Linda J. S. Allen, Armando Arciniega, and Priscilla E. Greenwood. Construction of equivalent stochastic differential equation models. *Stochastic Analysis and Applications*, 26(2):274–297, 2008.
- [2] Linda J.S. Allen. A primer on stochastic epidemic models: Formulation, numerical simulation, and analysis. *Infectious Disease Modelling*, 2(2):128–142, 2017. ISSN 2468-0427. doi: <https://doi.org/10.1016/j.idm.2017.03.001>.
- [3] Philip Bachman, R Devon Hjelm, and William Buchwalter. Learning representations by maximizing mutual information across views. *arXiv preprint arXiv:1906.00910*, 2019.
- [4] David Barber and Felix Agakov. The im algorithm: A variational approach to information maximization. In *Proceedings of the 16th International Conference on Neural Information Processing Systems*, NIPS’03, page 201–208, Cambridge, MA, USA, 2003. MIT Press.
- [5] Atilim Gunes Baydin, Barak A Pearlmutter, Alexey Andreyevich Radul, and Jeffrey Mark Siskind. Automatic differentiation in machine learning: a survey. *Journal of machine learning research*, 18, 2018.
- [6] Eli Bingham, Jonathan P Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul Szerlip, Paul Horsfall, and Noah D Goodman. Pyro: Deep universal probabilistic programming. *Journal of Machine Learning Research*, 2018.
- [7] Kathryn Chaloner and Isabella Verdinelli. Bayesian experimental design: A review. *Statistical Science*, pages 273–304, 1995.
- [8] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 1597–1607. PMLR, 13–18 Jul 2020.
- [9] Yanzhi Chen, Dinghuai Zhang, Michael U. Gutmann, Aaron Courville, and Zhanxing Zhu. Neural approximate sufficient statistics for implicit models. In *International Conference on Learning Representations*, 2021.
- [10] Alex R Cook, Gavin J Gibson, and Christopher A Gilligan. Optimal observation times in experimental epidemic processes. *Biometrics*, 64(3):860–868, 2008.
- [11] Kyle Cranmer, Johann Brehmer, and Gilles Louppe. The frontier of simulation-based inference. *Proceedings of the National Academy of Sciences*, 2020.
- [12] Mahasen B. Dehideniya, Christopher C. Drovandi, and James M. McGree. Optimal bayesian design for discriminating between models with intractable likelihoods in epidemiology. *Computational Statistics & Data Analysis*, 124:277–297, 2018. ISSN 0167-9473. doi: <https://doi.org/10.1016/j.csda.2018.03.004>.
- [13] Jacob Devlin, Ming Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, volume 1, pages 4171–4186, 2019.
- [14] Yann Dubois, Jonathan Gordon, and Andrew YK Foong. Neural Process Family. <http://yanndubs.github.io/Neural-Process-Family/>, September 2020.

- [15] Adam Foster, Martin Jankowiak, Elias Bingham, Paul Horsfall, Yee Whye Teh, Thomas Rainforth, and Noah Goodman. Variational Bayesian Optimal Experimental Design. In *Advances in Neural Information Processing Systems 32*, pages 14036–14047. Curran Associates, Inc., 2019.
- [16] Adam Foster, Martin Jankowiak, Matthew O’Meara, Yee Whye Teh, and Tom Rainforth. A unified stochastic gradient approach to designing bayesian-optimal experiments. In *International Conference on Artificial Intelligence and Statistics*, pages 2959–2969. PMLR, 2020.
- [17] Adam Foster, Desi R Ivanova, Ilyas Malik, and Tom Rainforth. Deep adaptive design: Amortizing sequential bayesian experimental design. *Proceedings of the 38th International Conference on Machine Learning (ICML)*, PMLR 139, 2021.
- [18] Andrew Gelman, John B Carlin, Hal S Stern, David B Dunson, Aki Vehtari, and Donald B Rubin. *Bayesian data analysis*. Chapman and Hall/CRC, 2013.
- [19] Matthew Graham and Amos Storkey. Asymptotically exact inference in differentiable generative models. In *Artificial Intelligence and Statistics*, pages 499–508. PMLR, 2017.
- [20] PeterJ. Green, Krzysztof Latuszynski, Marcelo Pereyra, and Christian P. Robert. Bayesian computation: a summary of the current state, and samples backwards and forwards. *Statistics and Computing*, 25(4):835–862, 2015. doi: 10.1007/s11222-015-9574-5.
- [21] Markus Hainy, Christopher C. Drovandi, and James M. McGree. Likelihood-free extensions for Bayesian sequentially designed experiments. In Joachim Kunert, Christine H. Müller, and Anthony C. Atkinson, editors, *mODa 11 - Advances in Model-Oriented Design and Analysis*, pages 153–161. Springer International Publishing, 2016.
- [22] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8): 1735–1780, 1997.
- [23] Xun Huan and Youssef Marzouk. Gradient-based stochastic optimization methods in bayesian experimental design. *International Journal for Uncertainty Quantification*, 4(6), 2014.
- [24] Xun Huan and Youssef M Marzouk. Simulation-based optimal bayesian experimental design for nonlinear systems. *Journal of Computational Physics*, 232(1):288–317, 2013.
- [25] Cheng Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Noam Shazeer, Ian Simon, Curtis Hawthorne, Andrew M. Dai, Matthew D. Hoffman, Monica Dinculescu, and Douglas Eck. Music transformer: Generating music with long-term structure, 2019. ISSN 23318422.
- [26] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [27] S. Kleinegesse and M.U. Gutmann. Efficient Bayesian experimental design for implicit models. In Kamalika Chaudhuri and Masashi Sugiyama, editors, *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 89 of *Proceedings of Machine Learning Research*, pages 1584–1592. PMLR, 2019.
- [28] Steven Kleinegesse and Michael Gutmann. Bayesian experimental design for implicit models by mutual information neural estimation. In *Proceedings of the 37th International Conference on Machine Learning*, Proceedings of Machine Learning Research, pages 5316–5326. PMLR, 2020.
- [29] Steven Kleinegesse and Michael U. Gutmann. Gradient-based bayesian experimental design for implicit models using mutual information lower bounds. *arXiv preprint arXiv:2105.04379*, 2021.
- [30] Steven Kleinegesse, Christopher Drovandi, and Michael U. Gutmann. Sequential Bayesian Experimental Design for Implicit Models via Mutual Information. *Bayesian Analysis*, pages 1 – 30, 2021. doi: 10.1214/20-BA1225.
- [31] Alexandre Lacoste, Alexandra Luccioni, Victor Schmidt, and Thomas Dandres. Quantifying the carbon emissions of machine learning. *arXiv preprint arXiv:1910.09700*, 2019.

- [32] Dennis V Lindley. On a measure of the information provided by an experiment. *The Annals of Mathematical Statistics*, pages 986–1005, 1956.
- [33] J. Lintusaari, M.U. Gutmann, R. Dutta, S. Kaski, and J. Corander. Fundamentals and recent developments in approximate Bayesian computation. *Systematic Biology*, 66(1):e66–e82, January 2017.
- [34] David McAllester and Karl Stratos. Formal limitations on the measurement of mutual information. In *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 875–884. PMLR, 2020.
- [35] Shakir Mohamed, Mihaela Rosca, Michael Figurnov, and Andriy Mnih. Monte carlo gradient estimation in machine learning. *Journal of Machine Learning Research*, 21(132):1–62, 2020.
- [36] Jay I Myung, Daniel R Cavagnaro, and Mark A Pitt. A tutorial on adaptive design optimization. *Journal of mathematical psychology*, 57(3-4):53–67, 2013.
- [37] Xuanlong Nguyen, Martin J. Wainwright, and Michael I. Jordan. Estimating divergence functionals and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information Theory*, 56(11), 2010. ISSN 00189448. doi: 10.1109/TIT.2010.2068870.
- [38] Antony Overstall and James McGree. Bayesian Design of Experiments for Intractable Likelihood Models Using Coupled Auxiliary Models and Multivariate Emulation. *Bayesian Analysis*, 15(1):103 – 131, 2020. doi: 10.1214/19-BA1144.
- [39] Emilio Parisotto, Francis Song, Jack Rae, Razvan Pascanu, Caglar Gulcehre, Siddhant Jayakumar, Max Jaderberg, Raphaël Lopez Kaufman, Aidan Clark, Seb Noury, Matthew Botvinick, Nicolas Heess, and Raia Hadsell. Stabilizing transformers for reinforcement learning. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 7487–7498. PMLR, 2020.
- [40] Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Noam Shazeer, Alexander Ku, and Dustin Tran. Image transformer. In *International Conference on Machine Learning*, pages 4055–4064. PMLR, 2018.
- [41] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems* 32, pages 8024–8035. Curran Associates, Inc., 2019.
- [42] Ben Poole, Sherjil Ozair, Aäron van den Oord, Alex Alemi, and George Tucker. On variational bounds of mutual information. In *International Conference on Machine Learning*, pages 5171–5180, 2019.
- [43] David J. Price, Nigel G. Bean, Joshua V. Ross, and Jonathan Tuke. On the efficient determination of optimal bayesian experimental designs using abc: A case study in optimal observation of epidemics. *Journal of Statistical Planning and Inference*, 172:1–15, May 2016.
- [44] David J Price, Nigel G Bean, Joshua V Ross, and Jonathan Tuke. An induced natural selection heuristic for finding optimal bayesian experimental designs. *Computational Statistics & Data Analysis*, 126:112–124, 2018.
- [45] Tom Rainforth. *Automating Inference, Learning, and Design using Probabilistic Programming*. PhD thesis, University of Oxford, 2017.
- [46] Tom Rainforth, Rob Cornish, Hongseok Yang, Andrew Warrington, and Frank Wood. On nesting monte carlo estimators. In *International Conference on Machine Learning*, pages 4267–4276. PMLR, 2018.



- [47] Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jon Shlens. Stand-alone self-attention in vision models. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [48] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *Proceedings of the 31st International Conference on Machine Learning*, volume 32, pages 1278–1286, 2014.
- [49] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- [50] Elizabeth G. Ryan, Christopher C. Drovandi, M. Helen Thompson, and Anthony N. Pettitt. Towards bayesian experimental design for nonlinear models that require a large number of sampling times. *Computational Statistics & Data Analysis*, 70:45–60, 2014. ISSN 0167-9473. doi: <https://doi.org/10.1016/j.csda.2013.08.017>.
- [51] Elizabeth G Ryan, Christopher C Drovandi, James M McGree, and Anthony N Pettitt. A review of modern computational algorithms for bayesian optimal design. *International Statistical Review*, 84(1):128–154, 2016.
- [52] Ben Shababo, Brooks Paige, Ari Pakman, and Liam Paninski. Bayesian inference and online experimental design for mapping neural microcircuits. In *Advances in Neural Information Processing Systems*, pages 1304–1312, 2013.
- [53] Xiaohong Sheng and Yu Hen Hu. Maximum likelihood multiple-source localization using acoustic energy measurements with wireless sensor networks. *IEEE Transactions on Signal Processing*, 2005. ISSN 1053587X. doi: 10.1109/TSP.2004.838930.
- [54] S.A. Sisson, Y. Fan, and M. Beaumont. *Handbook of Approximate Bayesian Computation*. Chapman & Hall/CRC Handbooks of Modern Statistical Methods. CRC Press, 2018. ISBN 9781351643467.
- [55] Jiaming Song and Stefano Ermon. Understanding the limitations of variational mutual information estimators. In *International Conference on Learning Representations*, 2020.
- [56] Owen Thomas, Ritabrata Dutta, Jukka Corander, Samuel Kaski, and Michael U Gutmann. Likelihood-free inference by ratio estimation. *arXiv preprint arXiv:1611.10242*, 2016.
- [57] Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- [58] Joep Vanlier, Christian A Tiemann, Peter AJ Hilbers, and Natal AW van Riel. A Bayesian approach to targeted experiment design. *Bioinformatics*, 28(8):1136–1142, 2012.
- [59] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [60] Benjamin T Vincent and Tom Rainforth. The DARC toolbox: automated, flexible, and efficient delayed and risky choice experiments using bayesian adaptive design. *PsyArXiv preprint*, 2017.
- [61] M. Zaharia, Andrew Chen, A. Davidson, A. Ghodsi, S. Hong, A. Konwinski, Siddharth Murching, Tomas Nykodym, Paul Ogilvie, Mani Parkhe, Fen Xie, and Corey Zumar. Accelerating the machine learning lifecycle with MLflow. *IEEE Data Eng. Bull.*, 41:39–45, 2018.
- [62] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabás Póczos, Ruslan Salakhutdinov, and Alexander J Smola. Deep sets. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS’17, 2017.
- [63] Jiabin Zhang, Sirui Bi, and Guannan Zhang. A stochastic approximate gradient ascent method for Bayesian experimental design with implicit models. In *The 24th International Conference on Artificial Intelligence and Statistics*, 2021.

## A Proofs

We present proof for all propositions made in the paper, restating each for convenience. We also include additional discussion on technical aspects of the paper.

### A.1 Unified objective for non-exchangeable experiments

**Proposition 1** (Generalized total expected information gain). *Consider the data generating distribution  $p(h_T|\theta, \pi) = \prod_{t=1:T} p(y_t|\theta, \xi_t, h_{t-1})$ , where  $\xi_t = \pi(h_{t-1})$  are the designs generated by the policy and, unlike in (4),  $y_t$  is allowed to depend on the history  $h_{t-1}$ . Then we can write (3) as*

$$\mathcal{I}_T(\pi) = \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} [\log p(h_T|\theta, \pi)] - \mathbb{E}_{p(h_T|\pi)} [\log p(h_T|\pi)]. \quad (6)$$

*Proof.* Starting with the definition of the total EIG (3) of a policy  $\pi$ :

$$\mathcal{I}_T(\pi) = \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} \left[ \sum_{t=1}^T I_{h_{t-1}}(\xi_t) \right] \quad (11)$$

we have by linearity of expectation

$$= \sum_{t=1}^T \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} [I_{h_{t-1}}(\xi_t)] \quad (12)$$

and since  $I_{h_{t-1}}$  doesn't depend on data acquired after  $t-1$  (the future doesn't influence the past)

$$= \sum_{t=1}^T \mathbb{E}_{p(\theta)p(h_{t-1}|\theta, \pi)} [I_{h_{t-1}}(\xi_t)] \quad (13)$$

which, applying Bayes rule, is equivalent to

$$= \sum_{t=1}^T \mathbb{E}_{p(h_{t-1}|\pi)p(\theta|h_{t-1})} [I_{h_{t-1}}(\xi_t)] \quad (14)$$

Next, using Bayes rule we similarly rearrange  $I_{h_{t-1}}$ :

$$I_{h_{t-1}}(\xi_t) = \mathbb{E}_{p(\theta|h_{t-1})p(y_t|\theta, \xi_t, h_{t-1})} \left[ \log \frac{p(y_t|\theta, \xi_t, h_{t-1})}{p(y_t|\xi_t, h_{t-1})} \right] \quad (15)$$

$$= \mathbb{E}_{p(\theta|h_{t-1})p(y_t|\theta, \xi_t, h_{t-1})} \left[ \log \frac{p(\theta|y_t, \xi_t, h_{t-1})}{p(\theta|h_{t-1})} \right] \quad (16)$$

$$= \mathbb{E}_{p(\theta|h_{t-1})p(y_t|\theta, \xi_t, h_{t-1})} [\log p(\theta|y_t, \xi_t, h_{t-1})] - \mathbb{E}_{p(\theta|h_{t-1})} [\log p(\theta|h_{t-1})] \quad (17)$$

$$= \mathbb{E}_{p(\theta|y_t, \xi_t, h_{t-1})p(y_t|\xi_t, h_{t-1})} [\log p(\theta|y_t, \xi_t, h_{t-1})] - \mathbb{E}_{p(\theta|h_{t-1})} [\log p(\theta|h_{t-1})] \quad (18)$$

and noting  $h_t = h_{t-1} \cup \{(\xi_t, y_t)\}$

$$= \mathbb{E}_{p(\theta|h_t)p(y_t|\xi_t, h_{t-1})} [\log p(\theta|h_t)] - \mathbb{E}_{p(\theta|h_{t-1})} [\log p(\theta|h_{t-1})] \quad (19)$$

$$= \mathbb{E}_{p(y_t|\xi_t, h_{t-1})} \left[ \mathbb{E}_{p(\theta|h_t)} [\log p(\theta|h_t)] - \mathbb{E}_{p(\theta|h_{t-1})} [\log p(\theta|h_{t-1})] \right] \quad (20)$$

Substituting this in (14), noting that  $\theta$  has already been integrated out, yields

$$\mathcal{I}_T(\pi) = \sum_{t=1}^T \mathbb{E}_{p(h_{t-1}|\pi)} \mathbb{E}_{p(y_t|\xi_t, h_{t-1})} \left[ \mathbb{E}_{p(\theta|h_t)} [\log p(\theta|h_t)] - \mathbb{E}_{p(\theta|h_{t-1})} [\log p(\theta|h_{t-1})] \right] \quad (21)$$

$$= \sum_{t=1}^T \mathbb{E}_{p(h_t|\pi)} \left[ \mathbb{E}_{p(\theta|h_t)} [\log p(\theta|h_t)] - \mathbb{E}_{p(\theta|h_{t-1})} [\log p(\theta|h_{t-1})] \right] \quad (22)$$

$$= \mathbb{E}_{p(h_T|\pi)} \left[ \sum_{t=1}^T \mathbb{E}_{p(\theta|h_t)} [\log p(\theta|h_t)] - \mathbb{E}_{p(\theta|h_{t-1})} [\log p(\theta|h_{t-1})] \right], \quad (23)$$

since we have a telescopic sum this simplifies to

$$= \mathbb{E}_{p(h_T|\pi)} \left[ \mathbb{E}_{p(\theta|h_T)} [\log p(\theta|h_T)] - \mathbb{E}_{p(\theta)} [\log p(\theta)] \right] \quad (24)$$

and finally we apply Bayes rule again to rewrite as

$$= \mathbb{E}_{p(h_T|\pi)p(\theta|h_T)} \left[ \log p(\theta|h_T) - \mathbb{E}_{p(\theta)} [\log p(\theta)] \right] \quad (25)$$

$$= \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} [\log p(\theta|h_T) - \log p(\theta)] \quad (26)$$

$$= \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} [\log p(h_T|\theta, \pi) - p(h_T|\pi)] \quad (27)$$

□

## A.2 Objective function as a mutual information

We provide some additional discussion on the interpretation of  $\mathcal{I}_T(\pi)$  in (6) as a mutual information. First,  $\mathcal{I}_T(\pi)$  is not a conventional mutual information between  $\theta$  and  $h_T$ . This is because, for the deterministic policy  $\pi$  considered in this paper, the random variable  $h_T$  does not have a density with respect to Lebesgue measure on  $\Xi^T \times \mathcal{Y}^T$ . Indeed, since the designs  $\xi_{1:T}$  are deterministic functions of the observations  $y_{1:T}$ , to express the sampling distribution of  $h_T$  we would have to use Dirac deltas, specifically

$$p(y_{1:T}, \xi_{1:T} | \theta, \pi) = \prod_{t=1}^T \delta_{\pi(h_{t-1})}(\xi_t) p(y_t | \theta, \xi_t, h_{t-1}). \quad (28)$$

Due to the presence of Dirac deltas, this is not a conventional probability density, and hence we do not regard  $\mathcal{I}_T(\pi)$  as the conventional mutual information between  $\theta$  and  $h_T$ .

We note that we *defined*  $p(h_T | \theta, \pi)$  in Proposition 1 differently to  $p(y_{1:T}, \xi_{1:T} | \theta, \pi)$  in (28). Specifically, our definition

$$p(h_T | \theta, \pi) = \prod_{t=1}^T p(y_t | \theta, \xi_t, h_{t-1}) \quad (29)$$

only involves probability densities for  $y_{1:T}$ , meaning that our  $p(h_T | \theta, \pi)$  is a well-defined probability density on  $\mathcal{Y}^T$ . Formally, we can treat the designs  $\xi_t$ , not as additional random variables, but as part of the density for  $y_{1:T}$ . Indeed, since the policy  $\pi$  is deterministic, it is possible to reconstruct  $h_{t-1}$  and  $\xi_t$  from  $y_{1:t-1}$  and  $\pi$ , so we could write  $p(y_t | \theta, y_{1:t-1}, \pi) := p(y_t | \theta, \xi_t, h_{t-1})$ . In this formulation, only  $y_{1:T}$  are regarded as random variables. This provides a formal justification for the form of  $p(h_T | \theta, \pi)$  that we give in Proposition 1. In this setting, we could formally identify  $\mathcal{I}_T(\pi)$  as the mutual information between  $\theta$  and  $y_{1:T}$ .

However, it is helpful to think of  $\mathcal{I}_T(\pi)$  as a mutual information between  $\theta$  and  $h_T$ , because this naturally leads to critics that have access to  $\theta$  and  $h_T$ , rather than  $\theta$  and  $y_{1:T}$ . This way of thinking also connects naturally to the case of stochastic policies, which we now discuss.

If we consider additional noise in the design process so that designs are no longer a deterministic function of past data, then  $\mathcal{I}_T(\pi)$  is the mutual information between  $\theta$  and  $h_T$ . In this case, we introduce an additional likelihood for designs  $p(\xi | \pi, h)$ , leading to the overall sampling distribution for the data

$$p(h_T | \theta, \pi) = \prod_{t=1}^T p(\xi_t | \pi, h_{t-1}) p(y_t | \theta, \xi_t, h_{t-1}). \quad (30)$$

Unlike in the deterministic case, this is valid probability density on  $\Xi^T \times \mathcal{Y}^T$ . If we now consider the mutual information between  $\theta$  and  $h_T$  for a fixed policy  $\pi$  we have

$$I(\theta, h_T) = \mathbb{E}_{p(\theta)p(h_T | \theta, \pi)} \left[ \log \frac{\prod_{t=1}^T p(\xi_t | \pi, h_{t-1}) p(y_t | \theta, \xi_t, h_{t-1})}{\int_{\Theta} p(\theta) \prod_{t=1}^T p(\xi_t | \pi, h_{t-1}) p(y_t | \theta, \xi_t, h_{t-1}) d\theta} \right] \quad (31)$$

$$= \mathbb{E}_{p(\theta)p(h_T | \theta, \pi)} \left[ \log \frac{\prod_{t=1}^T p(\xi_t | \pi, h_{t-1}) \prod_{t=1}^T p(y_t | \theta, \xi_t, h_{t-1})}{\prod_{t=1}^T p(\xi_t | \pi, h_{t-1}) \int_{\Theta} p(\theta) \prod_{t=1}^T p(y_t | \theta, \xi_t, h_{t-1}) d\theta} \right] \quad (32)$$

$$= \mathbb{E}_{p(\theta)p(h_T | \theta, \pi)} \left[ \log \frac{\prod_{t=1}^T p(y_t | \theta, \xi_t, h_{t-1})}{\int_{\Theta} p(\theta) \prod_{t=1}^T p(y_t | \theta, \xi_t, h_{t-1}) d\theta} \right] \quad (33)$$

noticing that the design likelihood terms cancel out in the integrand, and we reduce to the same integrand given in Proposition 1. Even when the policy is stochastic, the integrand in  $I(\theta, h_T)$  only involves terms of the form  $p(y_t | \theta, \xi_t, h_{t-1})$ , and the likelihood of the design process completely cancels. Thus, the stochasticity of the designs is only present in the sampling distribution  $p(h_T | \theta, \pi)$ . We therefore see that, as we consider the limiting case of  $p(\xi | \pi, h)$  as it approaches a deterministic policy, only the sampling distribution of designs in  $I(\theta, h_T)$  changes, with the integrand remaining the same. Under mild assumptions, then, the mutual information between  $\theta$  and  $h_T$  approaches  $\mathcal{I}_T(\pi)$  in this limit.

### A.3 NWJ and InfoNCE bounds

The next two propositions show that the two bounds—NWJ and InfoNCE—can be applied to the policy-based adaptive BOED setting.

**Proposition 2** (NWJ bound for implicit policy-based BOED). *For a design policy  $\pi$  and a critic function  $U : \mathcal{H}^T \times \Theta \rightarrow \mathbb{R}$ , let*

$$\mathcal{L}_T^{NWJ}(\pi, U) := \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} [U(h_T, \theta)] - e^{-1} \mathbb{E}_{p(\theta)p(h_T|\pi)} [\exp(U(h_T, \theta))], \quad (7)$$

*then  $\mathcal{I}_T(\pi) \geq \mathcal{L}_T^{NWJ}(\pi, U)$  holds for any  $U$ . Further, the inequality is tight for the optimal critic  $U_{NWJ}^*(h_T, \theta) = \log p(h_T|\theta, \pi) - \log p(h_T|\pi) + 1$ .*

*Proof.* Let  $\pi : \mathcal{H}^* \rightarrow \Xi$  be any (deterministic) policy taking histories  $h_t$  as inputs and returning a design  $\xi$  as output,  $U : \mathcal{H}^T \times \Theta \rightarrow \mathbb{R}$  be any function and define  $g(h_T, \theta) := \frac{\exp(U(h_T, \theta))}{\mathbb{E}_{p(h_T|\pi)} [\exp(U(h_T, \theta))]}$ .

First, we multiply the numerator and denominator of the unified objective (6) by  $g(h_T, \theta) > 0$

$$\mathcal{I}_T(\pi) = \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} \left[ \log \frac{p(h_T|\theta, \pi)}{p(h_T|\pi)} \right] \quad (34)$$

$$= \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} \log \left[ \frac{p(h_T|\theta, \pi)}{p(h_T|\pi)} \frac{g(h_T, \theta)}{g(h_T, \theta)} \right] \quad (35)$$

$$= \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} [\log g(h_T, \theta)] + \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} \left[ \log \frac{p(h_T|\theta, \pi)}{p(h_T|\pi)g(h_T, \theta)} \right] \quad (36)$$

Next, note that the second term is a KL divergence between two distributions

$$\mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} \left[ \log \frac{p(h_T|\theta, \pi)}{p(h_T|\pi)g(h_T, \theta)} \right] = \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} \left[ \log \frac{p(\theta)p(h_T|\theta, \pi)}{p(\theta)p(h_T|\pi)g(h_T, \theta)} \right] \quad (37)$$

$$= KL(p(\theta)p(h_T|\theta, \pi) || \hat{p}(h_T, \theta)) \geq 0 \quad (38)$$

where  $\hat{p}(h_T, \theta) = p(\theta)p(h_T|\pi)g(h_T, \theta)$  is a valid distribution since

$$\int p(\theta)p(h_T|\pi)g(h_T, \theta)d\theta dh_T = \mathbb{E}_{p(\theta)p(h_T|\pi)} \frac{\exp(U(h_T, \theta))}{\mathbb{E}_{p(h_T|\pi)} [\exp(U(h_T, \theta))]} \quad (39)$$

$$= \mathbb{E}_{p(\theta)} 1 = 1. \quad (40)$$

Therefore, we have

$$\mathcal{I}_T(\pi) \geq \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} [\log g(h_T, \theta)] \quad (41)$$

$$= \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} [U(h_T, \theta) - \log \mathbb{E}_{p(h_T|\pi)} \exp(U(h_T, \theta))] \quad (42)$$

$$= \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} [U(h_T, \theta)] - \mathbb{E}_{p(\theta)} [\log \mathbb{E}_{p(h_T|\pi)} \exp(U(h_T, \theta))] \quad (43)$$

Now using the inequality  $\log x \leq e^{-1}x$

$$\geq \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} [U(h_T, \theta)] - e^{-1} \mathbb{E}_{p(\theta)p(h_T|\pi)} [\exp(U(h_T, \theta))] \quad (44)$$

$$= \mathcal{L}_T^{NWJ}(\pi, U) \quad (45)$$

Finally, substituting  $U^*(h_T, \theta) = \log \frac{p(h_T|\theta, \pi)}{p(h_T|\pi)} + 1$  in the bound we get

$$\mathcal{L}_T^{NWJ}(\pi, U^*) = \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} \left[ \log \frac{p(h_T|\theta, \pi)}{p(h_T|\pi)} + 1 \right] - e^{-1} \mathbb{E}_{p(\theta)p(h_T|\pi)} \left[ \frac{p(h_T|\theta, \pi)}{p(h_T|\pi)} e^1 \right] \quad (46)$$

$$= \mathcal{I}_T(\pi) + 1 - \mathbb{E}_{p(\theta)p(h_T|\pi)} \left[ \frac{p(h_T|\theta, \pi)}{p(h_T|\pi)} \right] \quad (47)$$

$$= \mathcal{I}_T(\pi), \quad (48)$$

where we used  $\mathbb{E}_{p(\theta)p(h_T|\pi)} \left[ \frac{p(h_T|\theta, \pi)}{p(h_T|\pi)} \right] = \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} [1] = 1$ , establishing that the bound is tight for the optimal critic.  $\square$

**Proposition 3** (InfoNCE bound for implicit policy-based BOED). *Let  $\theta_{1:L} \sim p(\theta_{1:L}) = \prod_i p(\theta_i)$  be a set of contrastive samples where  $L \geq 1$ . For design policy  $\pi$  and critic function  $U: \mathcal{H}^T \times \Theta \rightarrow \mathbb{R}$ , let*

$$\mathcal{L}_T^{NCE}(\pi, U; L) := \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)} \mathbb{E}_{p(\theta_{1:L})} \left[ \log \frac{\exp(U(h_T, \theta_0))}{\frac{1}{L+1} \sum_{i=0}^L \exp(U(h_T, \theta_i))} \right], \quad (8)$$

*then  $\mathcal{I}_T(\pi) \geq \mathcal{L}_T^{NCE}(\pi, U; L)$  for any  $U$  and  $L \geq 1$ . Further, the optimal critic,  $U_{NCE}^*(h_T, \theta) = \log p(h_T|\theta, \pi) + c(h_T)$  where  $c(h_T)$  is any arbitrary function depending only on the history, recovers the sPCE bound in (5); the inequality is tight in the limit as  $L \rightarrow \infty$  for this optimal critic.*

*Proof.* Let  $\pi: \mathcal{H}^* \rightarrow \Xi$  be any (deterministic) policy taking histories  $h_t$  as inputs and returning a design  $\xi$  as output. Choose any function (critic)  $U: \mathcal{H}^T \times \Theta \rightarrow \mathbb{R}$ .

We introduce the shorthand

$$g(h_T, \theta_{0:L}) := \frac{\exp(U(h_T, \theta_0))}{\frac{1}{L+1} \sum_{i=0}^L \exp(U(h_T, \theta_i))} \quad (49)$$

Starting with the definition of the unified objective from Equation (6) we multiply its numerator and denominator by  $g(h_T, \theta_{0:L}) > 0$  to get

$$\mathcal{I}_T(\pi) = \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)} \left[ \log \frac{p(h_T|\theta_0, \pi)}{p(h_T|\pi)} \right] \quad (50)$$

where  $p(\theta_0)p(h_T|\theta_0, \pi) \equiv p(\theta)p(h_T|\theta, \pi)$

$$= \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)p(\theta_{1:L})} \left[ \log \frac{p(h_T|\theta_0, \pi)}{p(h_T|\pi)} \right] \quad (51)$$

$$= \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)p(\theta_{1:L})} \left[ \log \frac{p(h_T|\theta_0, \pi)g(h_T, \theta_{0:L})}{p(h_T|\pi)g(h_T, \theta_{0:L})} \right] \quad (52)$$

We next split the expectation into two terms one of which does not contain the unknown likelihoods and equals  $\mathcal{L}^{NCE}$

$$\begin{aligned} &= \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)p(\theta_{1:L})} \left[ \log \frac{p(h_T|\theta_0, \pi)}{p(h_T|\pi)g(h_T, \theta_{0:L})} \right] \\ &\quad + \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)p(\theta_{1:L})} [\log g(h_T, \theta_{0:L})] \\ &= \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)p(\theta_{1:L})} \left[ \log \frac{p(h_T|\theta_0, \pi)}{p(h_T|\pi)g(h_T, \theta_{0:L})} \right] + \mathcal{L}^{NCE}(\pi, U; L) \end{aligned} \quad (53)$$

We now show that the first term is a KL divergence and hence non-negative. To see why, first write

$$\mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)p(\theta_{1:L})} \left[ \log \frac{p(h_T|\theta_0, \pi)}{p(h_T|\pi)g(h_T, \theta_{0:L})} \right] \quad (54)$$

$$= \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)p(\theta_{1:L})} \left[ \log \frac{p(\theta_0)p(h_T|\theta_0, \pi)p(\theta_{1:L})}{p(\theta_0)p(h_T|\pi)p(\theta_{1:L})g(h_T, \theta_{0:L})} \right] \quad (55)$$

$$= \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)p(\theta_{1:L})} \left[ \log \frac{p(\theta_0)p(h_T|\theta_0, \pi)p(\theta_{1:L})}{\hat{p}(\theta_{0:L}, h_T|\pi)} \right] \quad (56)$$

$$= KL(p(h_T|\theta_0, \pi)p(\theta_{0:L}) || \hat{p}(\theta_{0:L}, h_T|\pi)). \quad (57)$$

and  $\hat{p}(\theta_{0:L}, h_T|\pi)$  is a valid distribution since

$$\int \hat{p}(\theta_{0:L}, h_T|\pi) d\theta_{0:L} dh_T = \int p(\theta_0)p(h_T|\pi)p(\theta_{1:L})g(h_T, \theta_{0:L}) d\theta_{0:L} dh_T \quad (58)$$

$$= \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)p(\theta_{1:L})} \left[ \frac{\exp(U(h_T, \theta_0))}{\frac{1}{L+1} \sum_{i=0}^L \exp(U(h_T, \theta_i))} \right], \quad (59)$$



because of the symmetry  $\theta_0 \stackrel{d}{=} \theta_j \forall j = 1, \dots, L$

$$= \frac{1}{L+1} \mathbb{E}_{p(\theta_0)p(h_T|\theta_0,\pi)p(\theta_{1:L})} \left[ \frac{\sum_{j=0}^L \exp(U(h_T, \theta_j))}{\frac{1}{L+1} \sum_{i=0}^L \exp(U(h_T, \theta_i))} \right] \quad (60)$$

$$= 1. \quad (61)$$

Thus we have established

$$\mathcal{I}_T(\pi) = KL(p(h_T|\theta_0, \pi)p(\theta_{0:L})||\hat{p}(\theta_{0:L}, h_T|\pi)) + \mathcal{L}_T^{NCE}(\pi, U; L) \geq \mathcal{L}_T^{NCE}(\pi, U; L). \quad (62)$$

Next, substituting  $U^*(h_T, \theta) = \log p(h_T|\theta, \pi) + c(h_T)$  in the definition of  $\mathcal{L}^{NCE}(\pi, U; L)$  we obtain

$$\mathcal{L}_T^{NCE}(\pi, U^*; L) = \mathbb{E}_{p(\theta_0)p(h_T|\theta_0,\pi)p(\theta_{1:L})} \left[ \frac{p(h_T|\theta_0, \pi) \exp(c(h_T))}{\frac{1}{L+1} \sum_{i=0}^L p(h_T|\theta_i, \pi) \exp(c(h_T))} \right] \quad (63)$$

$$= \mathbb{E}_{p(\theta_0)p(h_T|\theta_0,\pi)p(\theta_{1:L})} \left[ \frac{p(h_T|\theta_0, \pi)}{\frac{1}{L+1} \sum_{i=0}^L p(h_T|\theta_i, \pi)} \right], \quad (64)$$

which is exactly the sPCE bound (5), which is monotonically increasing in  $L$  and tight in the limit as  $L \rightarrow \infty$  [see 17, Theorem 2].  $\square$

#### A.4 A note on optimal critics

An interesting feature of our approach is that, for both the InfoNCE and NWJ bounds, the optimal critics do not depend on the policy. This is because we include the designs as explicit inputs to the critics. Indeed, we have

$$U_{\text{NCE}}^*(h_T, \theta) = \log \left( \prod_{t=1}^T p(y_t|\theta, \xi_t, h_{t-1}) \right) + c(h_T), \quad (65)$$

$$U_{\text{NWJ}}^*(h_T, \theta) = \log \left( \frac{\prod_{t=1}^T p(y_t|\theta, \xi_t, h_{t-1})}{\int_{\Theta} p(\theta) \prod_{t=1}^T p(y_t|\theta, \xi_t, h_{t-1}) d\theta} \right) + 1. \quad (66)$$

In previous work that utilized critics for gradient-based BOED [16, 28], it was typical to not treat the designs  $\xi_{1:T}$  as an input to the critic, which renders the optimal critic implicitly dependent on the designs. This makes more sense for static designs, for which the additional design input does not change. Our approach avoids an implicit dependence between policy and optimal critic which may be beneficial for the joint optimization.

## B Theoretical Comparison and Additional Bounds

Recently, a number of studies have discussed the challenges of estimating mutual information, in particular those associated with variational MI estimators [34, 42, 55].

Starting with the InfoNCE bound, it is trivial to show that the bound cannot exceed  $\log(L+1)$ , where  $L$  is the number of contrastive samples used to approximate the marginal in the denominator. Indeed,

$$\mathcal{L}_T^{NCE}(\pi, U; L) = \mathbb{E}_{p(\theta_0)p(h_T|\theta_0,\pi)} \mathbb{E}_{p(\theta_{1:L})} \left[ \log \frac{\exp(U(h_T, \theta_0))}{\frac{1}{L+1} \sum_{i=0}^L \exp(U(h_T, \theta_i))} \right] \quad (67)$$

$$\leq \log(L+1) + \mathbb{E}_{p(\theta_0)p(h_T|\theta_0,\pi)} \mathbb{E}_{p(\theta_{1:L})} \left[ \log \frac{\exp(U(h_T, \theta_0))}{\exp(U(h_T, \theta_0))} \right] \quad (68)$$

$$= \log(L+1) \quad (69)$$

This means that the corresponding Monte Carlo estimator will be highly biased whenever the true mutual information exceeds  $\log(L+1)$ , regardless of whether we have access to the optimal critic or not. This high bias estimator, however, comes with low variance [see e.g. 42, for discussion]. With

the optimal critic we would require exponential (in the MI) number of samples to accurately estimate the true mutual information.

It might appear at first that the NWJ bound might offer a better trade-off between bias and variance. Recall from the proof of Proposition 2, we have for the *optimal* critic

$$\mathcal{L}_T^{NWJ}(\pi, U^*) = \mathbb{E}_{p(\theta)p(h_T|\theta, \pi)} \left[ \log \frac{p(h_T|\theta, \pi)}{p(h_T|\pi)} \right] + 1 - e^{-1} \mathbb{E}_{p(\theta)p(h_T|\pi)} \left[ \frac{p(h_T|\theta, \pi)}{p(h_T|\pi)} e^1 \right], \quad (70)$$

of which we form a Monte carlo estimate using  $N$  ( $M$ ) samples for the first (second) term, respectively

$$\approx \frac{1}{N} \sum_{n=1}^N \log \frac{p(h_{T,n}|\theta_n, \pi)}{p(h_{T,n}|\pi)} + \left( 1 - \frac{1}{M} \sum_{m=1}^M \log \frac{p(h_{T,m}|\theta_m, \pi)}{p(h_{T,m}|\pi)} \right), \quad (71)$$

where  $\theta_n, h_{T,n} \sim p(\theta)p(h_T|\theta, \pi)$  are samples from the joint distribution and  $\theta_m, h_{T,m} \sim p(\theta)p(h_T|\pi)$  are samples from the product of marginals. The first term is a Monte Carlo estimate of the mutual information, while the second has mean zero, meaning that this estimator is unbiased. The second term, however has variance which grows exponentially with the value of the (true) mutual information [see Theorem 2 in 55]. What this means is that even with an optimal critic, we will need an exponential (in the MI) number of samples to control the variance of the NWJ estimator. One might then hope that the variance can be reduced when using a sub-optimal critic at the cost of introducing some (hopefully small) bias. Unfortunately, according to a recent result [see Theorems 3.1 and 4.1 in 34, and the discussion therein], it is not possible to guarantee that a likelihood-free lower bound on the mutual information can exceed  $\log(N)$ . Indeed, the authors show theoretically and empirically that all high-confidence distribution-free lower bounds on the mutual information require exponential (in the the MI) number of samples.

Constructing a better lower bound on the mutual information—one that does not need exponential number of samples—therefore, requires us to make additional assumptions. Foster et al. [17] propose one such bound, namely the sequential Adaptive Constrative Estimation (sACE). The sACE bound introduces a proposal distribution  $q(\theta; h_T)$ , which aims to approximate the posterior  $p(\theta|h_T)$ . Since implicit models were not the focus of the work in [17] the proposed bound, relies on analytically available likelihood. The following proposition shows we can derive a likelihood-free version of the sACE bound.

**Proposition 4** (Sequential Likelihood-free ACE). *For a design function  $\pi$ , a critic function  $U$ , a number of contrastive samples  $L \geq 1$ , and a proposal  $q(\theta; h_T)$ , we have the sequential Likelihood-free Adaptive Contrastive Estimation (sLACE) lower bound*

$$\mathcal{L}_T^{sLACE}(\pi, U, q; L) := \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)q(\theta_{1:L}; h_T)} \left[ \log \frac{U(h_T, \theta_0)}{\frac{1}{L+1} \sum_{\ell=0}^L \frac{U(h_T, \theta_\ell)p(\theta_\ell)}{q(\theta_\ell; h_T)}} \right] \leq \mathcal{I}_T(\pi). \quad (72)$$

The bound is tight as  $L \rightarrow \infty$  for the optimal critic  $U^*(h_T, \theta) = \log p(h_T|\theta, \pi) + c(h_T)$ , where  $c(h_T)$  is arbitrary. In addition, if  $q(\theta; h_T) = p(\theta|h_T)$ , the bound is tight for the optimal critic  $U^*(h_T, \theta)$  with any  $L \geq 0$ .

*Proof.* The proof follows similar arguments to the ones in Propositions 2 and 3. First let

$$g(h_T, \theta_{0:L}) := \frac{U(h_T, \theta_0)}{\frac{1}{L+1} \sum_{\ell=0}^L \frac{U(h_T, \theta_\ell)p(\theta_\ell)}{q(\theta_\ell; h_T)}} \quad (73)$$

Starting with the definition of the EIG:

$$\mathcal{I}_T(\pi) = \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)} \left[ \log \frac{p(h_T|\theta_0, \pi)}{p(h_T|\pi)} \right] \quad (74)$$

since  $q(\theta_i; h_T)$  is a valid density

$$= \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)q(\theta_{1:L}; h_T)} \left[ \log \frac{p(h_T|\theta_0, \pi)}{p(h_T|\pi)} \right] \quad (75)$$

multiplying its numerator and denominator inside the log by  $g(h_T, \theta_{0:L}) > 0$

$$= \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)q(\theta_{1:L}; h_T)} \left[ \log \frac{p(h_T|\theta_0, \pi)g(h_T, \theta_{0:L})}{p(h_T|\pi)g(h_T, \theta_{0:L})} \right] \quad (76)$$

$$= \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)q(\theta_{1:L}; h_T)} [\log g(h_T, \theta_{0:L})] \\ + \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)q(\theta_{1:L}; h_T)} \left[ \log \frac{p(h_T|\theta_0, \pi)}{p(h_T|\pi)g(h_T, \theta_{0:L})} \right] \quad (77)$$

The first term is exactly the sLACE bound,  $\mathcal{L}_T^{\text{sLACE}}(\pi, U, q; L)$ . We now show that the second term is a KL divergence between two distributions and hence non-negative. To see this

$$\mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)q(\theta_{1:L}; h_T)} \left[ \log \frac{p(h_T|\theta_0, \pi)}{p(h_T|\pi)g(h_T, \theta_{0:L})} \right] \quad (78)$$

$$= \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)q(\theta_{1:L}; h_T)} \left[ \log \frac{p(\theta_0)p(h_T|\theta_0, \pi)q(\theta_{1:L}; h_T)}{p(h_T|\pi)g(h_T, \theta_{0:L})p(\theta_0)q(\theta_{1:L}; h_T)} \right] \quad (79)$$

$$= KL(p(\theta_0)p(h_T|\theta_0, \pi)q(\theta_{1:L}; h_T) || \hat{p}(h_T, \theta_{0:L})), \quad (80)$$

since  $\hat{p}(h_T, \theta_{0:L}) := p(h_T|\pi)g(h_T, \theta_{0:L})p(\theta_0)q(\theta_{1:L}; h_T)$  is a valid density. Indeed:

$$\int \hat{p}(h_T, \theta_{0:L}) dh_T d\theta_{0:L} = \mathbb{E}_{q(\theta_{1:L}; h_T)p(h_T|\pi)} [p(\theta_0)g(h_T, \theta_{0:L})] \quad (81)$$

$$= \mathbb{E}_{q(\theta_{1:L}; h_T)p(h_T|\pi)} \left[ p(\theta_0) \frac{U(h_T, \theta_0)}{\frac{1}{L+1} \sum_{\ell=0}^L \frac{U(h_T, \theta_\ell)p(\theta_\ell)}{q(\theta_\ell; h_T)}} \right] \quad (82)$$

$$= \mathbb{E}_{q(\theta_{0:L}; h_T)p(h_T|\pi)} \left[ \frac{\frac{U(h_T, \theta_0)p(\theta_0)}{q(\theta_0; h_T)}}{\frac{1}{L+1} \sum_{\ell=0}^L \frac{U(h_T, \theta_\ell)p(\theta_\ell)}{q(\theta_\ell; h_T)}} \right] \quad (83)$$

by symmetry

$$= \mathbb{E}_{q(\theta_{0:L}; h_T)p(h_T|\pi)} \left[ \frac{\frac{1}{L+1} \sum_{\ell=0}^L \frac{U(h_T, \theta_\ell)p(\theta_\ell)}{q(\theta_\ell; h_T)}}{\frac{1}{L+1} \sum_{\ell=0}^L \frac{U(h_T, \theta_\ell)p(\theta_\ell)}{q(\theta_\ell; h_T)}} \right] \quad (84)$$

$$= 1. \quad (85)$$

With the optimal critic we recover the sACE bound from [17], which under mild conditions converges to the mutual information  $\mathcal{I}_T(\pi)$ . To see that start by writing

$$\mathcal{L}_T^{\text{sLACE}}(\pi, U^*, q; L) = \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)q(\theta_{1:L}; h_T)} \left[ \log \frac{p(h_T|\theta_0, \pi)}{\frac{1}{L+1} \sum_{\ell=0}^L \frac{p(h_T|\theta_\ell, \pi)p(\theta_\ell)}{q(\theta_\ell; h_T)}} \right]. \quad (86)$$

The denominator is a consistent estimator of the marginal, provided that each term in the sum is bounded, and so by the Strong Law of Large Numbers we have

$$\frac{1}{L+1} \sum_{\ell=0}^L \frac{p(h_T|\theta_\ell, \pi)p(\theta_\ell)}{q(\theta_\ell; h_T)} \rightarrow p(h_T|\pi) \text{ a.s. as } L \rightarrow \infty, \quad (87)$$

which establishes point-wise convergence of the integrand to  $p(h_T|\theta_0, \pi)/p(h_T|\pi)$ . We can apply Bounded convergence theorem to establish  $\mathcal{L}_T^{\text{sACE}}(\pi, U^*, q; L) \rightarrow \mathcal{I}_T(\pi)$  as  $L \rightarrow \infty$ .

If in addition  $q(\theta; h_T) = p(\theta|h_T)$  we have by Bayes rule:

$$\mathcal{L}_T^{\text{sLACE}}(\pi, U^*, q; L) = \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)p(\theta_{1:L}|h_T)} \left[ \log \frac{p(h_T|\theta_0, \pi)}{\frac{1}{L+1} \sum_{\ell=0}^L \frac{p(h_T|\theta_\ell, \pi)p(\theta_\ell)}{p(\theta_\ell|h_T)}} \right] \quad (88)$$

$$= \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)p(\theta_{1:L}|h_T)} \left[ \log \frac{p(h_T|\theta_0, \pi)}{\frac{1}{L+1} \sum_{\ell=0}^L p(h_T|\pi)} \right] \quad (89)$$

$$= \mathbb{E}_{p(\theta_0)p(h_T|\theta_0, \pi)} \left[ \log \frac{p(h_T|\theta_0, \pi)}{p(h_T|\pi)} \right] \quad (90)$$

$$= \mathcal{I}_T(\pi) \quad \forall L \geq 0. \quad (91)$$

□

In practice, we parameterize the policy, the critic and the density of the proposal distribution by neural networks  $\pi_\phi$ ,  $U_\psi$  and  $q_\zeta$  and optimize  $\mathcal{L}_T^{\text{sLACE}}$  with respect to the parameters of these networks,  $\phi$ ,  $\psi$  and  $\zeta$  with SGA. As before, optimizing with respect to  $\phi$  improves the quality of the designs, proposed by the policy, whilst optimizing with respect to  $\psi$  and  $\zeta$  tightens the bound. If the parametric density  $q_\zeta$  and the critic  $U_\psi$  are expressive enough, so that we can recover the optimal critic and the true posterior, then the bound is tight for any number of contrastive samples  $L$ . If, on the other hand, we fix  $q_\zeta(\theta; h_T) = p(\theta)$  instead of training it, then we recover the InfoNCE bound. Therefore, as long as  $q_\zeta$  approximates the posterior better than the prior, then even an imperfect proposal  $q_\zeta$  can benefit training.

In addition to introducing another set of optimizable parameters,  $\zeta$ , the sLACE bound assumes that we know the prior  $p(\theta)$  and can evaluate its density.

## C Neural architecture

### C.1 Permutation invariance of the critic for exchangeable experiments

We show that if the BOED problem is exchangeable then the critic function  $U$  should be permutation-invariant.

**Proposition 5** (Permutation invariance). *Let  $\sigma$  be a permutation acting on a history  $h_T^1$  yielding  $h_T^2 = \{(\xi_{\sigma(i)}, y_{\sigma(i)})\}_{i=1}^T$ . If the data generating process is conditionally independent of its past given  $\theta$ , then the optimal critics for both (7) and (8) are invariant under permutations of the history, i.e.*

$$p(\theta) \prod_{t=1}^T p(y_t | \theta, \xi_t(h_{t-1}), h_{t-1}) = p(\theta) \prod_{t=1}^T p(y_t | \theta, \xi_t) \implies U^*(h_T^1, \theta) = U^*(h_T^2, \theta). \quad (92)$$

*Proof.* This is a direct consequence from the form of the optimal critics. To see this formally, let  $h_T^1$  be a history and  $h_T^2$  be a permutation of it.

Starting with the InfoNCE bound we have

$$U_{\text{NCE}}^*(h_T^1, \theta) = \log p(h_T^1 | \theta, \pi) + c(h_T^1) \quad (93)$$

$$= \log \prod_{t=1}^T p(y_t | \theta, \xi_t) + c(\{(\xi_t, y_t)\}_{t=1}^T) \quad (94)$$

since  $c(h_T)$  is arbitrary, we can choose it to be permutation invariant

$$= \log \prod_{t=1}^T p(y_{\sigma(t)} | \theta, \xi_{\sigma(t)}) + c(\{(\xi_{\sigma(t)}, y_{\sigma(t)})\}_{t=1}^T) \quad (95)$$

$$= \log p(h_T^2 | \theta, \pi) + c(h_T^2) \quad (96)$$

$$= U_{\text{NCE}}^*(h_T^2, \theta) \quad (97)$$

Similarly, for the optimal critic of the NWJ bound we have

$$U_{\text{NWJ}}^*(h_T^1, \theta) = \log \frac{p(h_T^1 | \theta, \pi)}{p(h_T^1 | \pi)} + 1 \quad (98)$$

$$= \log \frac{\prod_{t=1}^T p(y_t | \theta, \xi_t)}{\mathbb{E}_{p(\theta)} \left[ \prod_{s=1}^T p(y_s | \theta, \xi_s) \right]} + 1 \quad (99)$$

$$= \log \frac{\prod_{t=1}^T p(y_{\sigma(t)} | \theta, \xi_{\sigma(t)})}{\mathbb{E}_{p(\theta)} \left[ \prod_{s=1}^T p(y_{\sigma(s)} | \theta, \xi_{\sigma(s)}) \right]} + 1 \quad (100)$$

$$= \log \frac{p(h_T^2 | \theta, \pi)}{p(h_T^2 | \pi)} + 1 = U_{\text{NWJ}}^*(h_T^2, \theta). \quad (101)$$

□

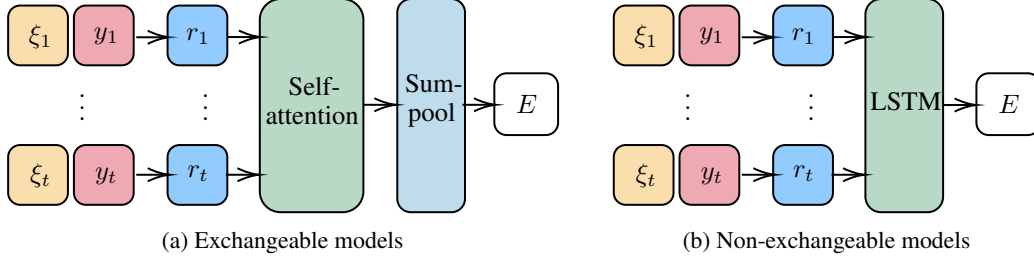


Figure 6: History encoder architectures for different classes of models. When conditional independence of the experiments holds, we use self-attention, followed by sum-pooling, making the history encoder permutation invariant. When experiments are not conditionally independent we use LSTM and only keep its last hidden state. We train two separate history encoders—one for the design network  $\pi_\phi$  and one for the critic network  $U_\psi$ , although we note that all the weights except those in the head layers can be shared.

To the best of our knowledge, we are the first to propose a critic architecture that is tailored to BOED problems with exchangeable models. Previous work in the static BOED setting, where MI information objective is optimized with variational lower bounds and thus require the training of critics [e.g. 28, 63], did not discuss what an appropriate critic architecture might be. In particular, in all experiments [28, 63] use a generic architecture for both exchangeable and non-exchangeable problems. An expressive enough generic architecture should be able to obtain the optimal critic, and thus achieve a tight bound, however, the optimisation process will be considerably more difficult as the network needs to learn this key invariance structure. We therefore recommend using permutation invariant architectures whenever the model is exchangeable, especially if achieving tight bounds (and therefore learning an optimal critic) is of importance.

## C.2 Further details on the history encoder

Figure 6 shows the history encoders we use in the policy network  $\pi_\phi$  and the critic network  $U_\psi$ . First, we encode the individual design-outcome pairs,  $(\xi_t, y_t)$ , with an MLP, which gives us a vector of representations  $r_t \in \mathbb{R}^m$ , where  $m$  is the encoding dimension we have selected. The representations  $\{r_i\}_{i=1}^t$  are row-stacked into a matrix  $R$  of dimension  $t \times m$ , which we then aggregate back to a vector of size  $m$  by an appropriate layer(s).

When conditional independence of the experiments holds, we apply 8-head self-attention, based on the Image Transformer [40] and as implemented by [14]. Applying self-attention leaves the dimension of the matrix  $R$  unchanged. We then apply sum-pooling across time  $t$ , which gives us the final encoding vector  $E \in \mathbb{R}^m$ .

When experiments are not conditionally independent, we pass the matrix  $R$  through an LSTM with two hidden layers and hidden state of size  $m$  (see the **LSTM module in Pytorch** for more details). The LSTM returns hidden state vectors associated with the history  $h_t$  for each  $t$ ; we keep the last hidden state of the last layer, which is our final encoding vector  $E \in \mathbb{R}^m$ .

In both cases the resulting encoding  $E$  is a vector of size  $m$ . It is passed through final fully connected "head" layers, which output either a design (in the case of the policy) or a vector (in the case of the critic). We train two separate history encoders—one for the design network  $\pi_\phi$  and one for the critic network  $U_\psi$ , although we note that all the weights except those in the head layers can be shared.

## D Experiments

### D.1 Computational resources

All of the experiments were implemented in Python using open-source software. All estimators and models were implemented in PyTorch [41] (BSD license) and Pyro [6] (Apache License Version 2.0), whilst MIFlow [61] (Apache License Version 2.0) was used for experiment tracking and management. The self-attention architecture from [14] was used to implement the self-attention mechanisms in the



design and critic networks. For full details on package versions, environment set-up and commands for running the code, see instructions in the `README.md` file.

Experiments were ran on internal GPU clusters, consisting of GeForce RTX 3090 (24GB memory), GeForce RTX 2080 Ti (11GB memory) and GeForce GTX 1080 Ti GPUs (11GB memory).

The deployment-time of iDAD (Table 4) was estimated on a lightweight CPU machine with the following specifications

Processor	2.8 GHz Quad-Core Intel Core i7
Memory	16 GB
Operating system	macOS Big Sur v11.2.3

## D.2 CO2 Emission Related to Experiments

Experiments were conducted using a private infrastructure, which has an estimated carbon efficiency of 0.432 kgCO<sub>2</sub>eq/kWh. A cumulative of 160 hours of computation was performed on hardware of type RTX 2080 Ti (TDP of 250W), or similar. The training time of each experiment (including the baselines that require optimization), took on average between 1-3 GPU hours, depending on the number of experiments  $T$ .

Total emissions are estimated to be 17.28 kgCO<sub>2</sub>eq of which 0% was directly offset.

Estimations were conducted using the [Machine Learning Impact calculator](#) presented in [31].

## D.3 Traditional sequential BOED with variational posterior estimator

The variational posterior estimator from [15] is based on the Barbar-Agakov lower bound [4], which takes the form

$$\mathcal{L}^{\text{post}}(\xi, q_\psi) = \mathbb{E}_{p(\theta)p(y|\theta, \xi)} \left[ \log \frac{q_\psi(\theta; y, \xi)}{p(\theta)} \right] \leq \mathcal{I}(\xi), \quad (102)$$

where  $q_\psi$  is any normalized distribution over the parameters  $\theta$ . The bound is tight when  $q_\psi(\theta; y, \xi) = p(\theta|y, \xi)$ , i.e. if we can recover the true posterior. We assume mean-field variational family and optimize the parameters  $\psi$  by maximizing the bound (102) using stochastic gradient schemes. Simultaneously we optimize the bound with respect to the design variable  $\xi$  to select the optimal design  $\xi^*$ . At the inference stage, denoting by  $y^*$  the outcome of experiment  $\xi^*$ , we obtain an approximate posterior by evaluating  $q_\psi(\theta; y^*, \xi^*)$ , i.e. we reuse the learnt variational posterior. We repeat this process at each stage of the experiments by substituting the the approximate posterior,  $q_\psi(\theta; y^*, \xi^*)$ , as the prior in (102).

## D.4 Location Finding

In this experiment we have  $K$  hidden objects (*sources*) in  $\mathbb{R}^2$  and we wish to learn their locations,  $\theta = \{\theta_1, \dots, \theta_K\}$ . The number of sources,  $K$ , is assumed to be known. Each source emits a signal with intensity obeying the inverse-square law. Put differently, if a source is located at  $\theta_k$  and we perform a measurement at a point  $\xi$ , the signal strength emitted from that source only will be proportional to  $\frac{1}{\|\theta_k - \xi\|^2}$ . The total intensity at location  $\xi$ , emitted from all  $K$  sources, is a superposition of the individual ones

$$\mu(\theta, \xi) = b + \sum_{k=1}^K \frac{\alpha_k}{m + \|\theta_k - \xi\|^2}, \quad (103)$$

where  $\alpha_k$  can be known constants or random variables,  $b > 0$  is a constant background signal and  $m$  is a constant, controlling the maximum signal.

We place a standard normal prior on each of the location parameters  $\theta_k$  and we observe the log-total intensity with some Gaussian noise. We therefore have the following prior and likelihood:

$$\theta_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0_d, I_d) \quad \log y \mid \theta, \quad \xi \sim \mathcal{N}(\log \mu(\theta, \xi), \sigma^2) \quad (104)$$

#### D.4.1 Training details

All our experiments are performed with the following model hyperparameters

Parameter	Value
Number of sources, $K$	2
$\alpha_k$	$1 \forall k$
Max signal, $m$	$10^{-4}$
Base signal, $b$	$10^{-1}$
Observation noise scale, $\sigma$	0.5

The architecture of the design network  $\pi_\phi$  used in Table 2 and all its hyperparameters are in the following tables. For the encoder of the design-outcome pairs we used the following:

Layer	Description	iDAD, InfoNCE	iDAD, NWJ	Activation
Input	$\xi, y$	3	3	-
H1	Fully connected	64	64	ReLU
H2	Fully connected	512	512	ReLU
Output	Fully connected	64	64	-
Attention	8 heads	64	64	-

The output of the encoder,  $R(h_t)$ , is fed into an emitter network, for which we used the following:

Layer	Description	iDAD, InfoNCE	iDAD, NWJ	Activation
Input	$R(h_t)$	64	64	-
H1	Fully connected	256	256	ReLU
H2	Fully connected	64	64	ReLU
Output	Fully connected	2	2	-

The architecture of the critic network  $U_\psi$  used in Table 2 and all its hyperparameters are in the tables that follow. First, the encoder network of the latent variables is:

Layer	Description	iDAD, InfoNCE	iDAD, NWJ	Activation
Input	$\theta$	4	4	-
H1	Fully connected	16	16	ReLU
H2	Fully connected	64	64	ReLU
H3	Fully connected	512	512	ReLU
Output	Fully connected	64	64	-

For the design-outcome pairs encoder we use the same architecture as in the design network, namely:

Layer	Description	iDAD, InfoNCE	iDAD, NWJ	Activation
Input	$\xi, y$	3	3	-
H1	Fully connected	64	64	ReLU
H2	Fully connected	512	512	ReLU
Output	Fully connected	64	64	-
Attention	8 heads	64	64	-

The output of the encoder,  $R(h_t)$ , is fed into fully connected head layers:

Layer	Description	iDAD, InfoNCE	iDAD, NWJ	Activation
Input	$R(h_t)$	64	64	-
H1	Fully connected	1024	1024	ReLU
H2	Fully connected	512	512	ReLU
H3	Fully connected	512	512	ReLU
Output	Fully connected	64	64	-

The optimisation was performed with Adam [26] with ReduceLROnPlateau learning rate scheduler, with the following hyperparameters:

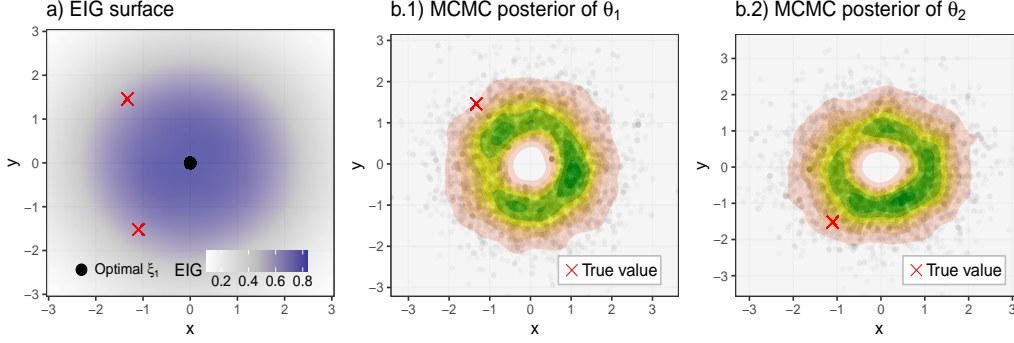


Figure 7: a): EIG surface induced by the prior; b) Samples from  $p(\theta|\xi_1, y_1)$ —the posterior distribution of the locations, after performing experiment  $\xi_1$  and observing  $y_1$ , along with a KDE.

Parameter	iDAD, InfoNCE	iDAD, NWJ
Batch size	2048	2048
Number of contrastive/negative samples	2047	2047
Number of gradient steps	100000	100000
Initial learning rate (LR)	0.0005	0.0005
LR annealing factor	0.8	0.8
LR annealing frequency (if no improvement)	2000	2000

#### D.4.2 Performance of the variational baseline

As we saw in Table 2, this variational approach to (myopic) adaptive BOED performed very poorly, despite its large computational budget. The likely reason for that is that the mean-field variational approximation cannot adequately capture the complex non-Gaussian posterior of this problem. Figure 7 clearly demonstrates this: before any data is observed it is optimal to sample at the origin (since the prior is centered at it). After observing a low signal (the locations in this example are not close to the origin), we can only conclude that the sources are not within a small radius of the origin, but anywhere outside of it would be a plausible location, as indeed indicated by the fitted posteriors.

#### D.4.3 Hyperparameter selection

We did not perform extensive hyperparameters search; in particular, the network sizes were guided by two hyperparameters: hidden-dimension ( $HD = 512$ ) and encoding dimension ( $ED = 64$ ). We set-up all the networks to scale up with the number of experiments as follows:

- Design-outcome encoder has three hidden layers of sizes  $[64, HD, ED]$ .
- Design emitter network has three hidden layers of sizes  $[HD/2, ED, 2]$ , where 2 is the dimension of the design variable.
- The latent encoder for the critic network has four hidden layers of sizes  $[16, 64, HD, ED]$ .
- The critic design-outcome encoder’s head layer has four hidden layers of sizes  $[HD \times \log(T), HD \times \log(T)/2, HD, ED]$ .

Since our multi-head attention layer has 8 heads, the encoding dimension we use has to be a multiple of 8. In addition to  $ED = 64$  we tried  $ED = 32$  which provided marginally worse results. We did not try other values for these hyperparameters.

For the learning rate, we tried 0.001, which was too high, as well as 0.0005 (which we selected) and 0.0001 (which yielded very similar results).

We performed similar level of hyperparameter tuning for all trainable baselines as well (DAD, MINEBED and SG-BOED).

Table 6: Upper bounds on the total information,  $\mathcal{I}_{10}(\pi)$ , for the location finding experiment in Section 5.1. The bounds were estimated using  $L = 5 \times 10^5$  contrastive samples. Errors indicate  $\pm 1$  s.e. estimated over 4096 histories (128 for variational). Lower bounds are presented in Table 2.

Method \ $\theta$ dim.	4D	6D	10D	20D
Random	$4.794 \pm 0.041$	$3.506 \pm 0.004$	$1.895 \pm 0.003$	$0.552 \pm 0.001$
MINEBED	$5.522 \pm 0.028$	$4.229 \pm 0.029$	$2.459 \pm 0.029$	$0.801 \pm 0.019$
SG-BOED	$5.549 \pm 0.028$	$4.220 \pm 0.030$	$2.455 \pm 0.029$	$0.803 \pm 0.019$
Variational	$4.644 \pm 0.146$	$3.626 \pm 0.167$	$2.181 \pm 0.152$	$0.669 \pm 0.097$
iDAD (NWJ)	$7.806 \pm 0.050$	$5.851 \pm 0.041$	<b><math>3.264 \pm 0.039</math></b>	$0.877 \pm 0.022$
iDAD (InfoNCE)	<b><math>7.863 \pm 0.043</math></b>	<b><math>6.068 \pm 0.039</math></b>	<b><math>3.257 \pm 0.040</math></b>	<b><math>0.872 \pm 0.020</math></b>
DAD	$8.034 \pm 0.038$	$6.310 \pm 0.031$	$3.358 \pm 0.040$	$0.953 \pm 0.022$

Table 7: Upper and lower bounds on the total information,  $\mathcal{I}_{20}(\pi)$ , for the location finding experiment in 2D from Section 5.1. The bounds were estimated using  $L = 5 \times 10^5$  contrastive samples. Errors indicate  $\pm 1$  s.e. estimated over 4096 histories.

Method	Lower bound	Upper bound
Random	$7.000 \pm 0.034$	$7.020 \pm 0.034$
MINEBED	$7.672 \pm 0.030$	$7.690 \pm 0.031$
SG-BOED	$7.701 \pm 0.030$	$7.728 \pm 0.031$
iDAD (NWJ)	<b><math>9.961 \pm 0.033</math></b>	<b><math>10.372 \pm 0.048</math></b>
iDAD (InfoNCE)	<b><math>10.075 \pm 0.032</math></b>	<b><math>10.463 \pm 0.043</math></b>
DAD	$10.424 \pm 0.031$	$10.996 \pm 0.049$

#### D.4.4 Further ablation studies

**Scalability with number of experiments.** We first demonstrate that iDAD can scale to a larger number of experiments  $T$ . We train policy networks to perform  $T = 20$  experiments and compare them to baselines in Table 7. We omit the variational baseline as it is too computationally costly to run for a large enough number of histories, and as we saw in the previous subsection, it is not particularly suited to this model.

**Training stability.** To assess the robustness of the results and the stability of the training process, we trained 5 additional iDAD networks with each of the two bounds, using different seeds but the same hyperparameters (described in Subsection D.4.1) we used to produce the results of the location finding experiment in 2D (Table 2 in the main text). We report upper and lower bounds on the mutual information along with their mean and standard error in the table below.

Estimator	Bound	Run 1	Run 2	Run 3	Run 4	Run 5	Mean	SE
InfoNCE	Lower	7.826	7.682	7.856	7.713	7.804	<b>7.776</b>	<b>0.034</b>
InfoNCE	Upper	7.933	7.791	7.856	7.807	7.925	<b>7.862</b>	<b>0.029</b>
NWJ	Lower	7.820	7.545	7.592	7.555	7.691	<b>7.641</b>	<b>0.052</b>
NWJ	Upper	7.976	7.640	7.669	7.651	7.800	<b>7.747</b>	<b>0.064</b>

We can see that the iDAD networks trained with InfoNCE are highly stable, with the 5 additional runs achieving very similar mutual information values to each other and to the iDAD network used the report the results in the main paper. The performance of the iDAD networks trained with the NWJ bound is more variable and empirically achieve slightly lower average value of mutual information. This higher variance is in-line with the discussion in Section B.

We similarly verify the robustness of the static baselines, reporting the results in the table below:

Table 8: Ablation study on the performance of iDAD as a function of training time for the location finding experiment.

Training budget	MI lower bound
0.1%	3.38
1.0%	6.09
2.0%	6.46
4.0%	6.81
8.0%	7.08
16.0%	7.33
32.0%	7.56
64.0%	7.78
100.0%	7.82

Estimator	Bound	Run 1	Run 2	Run 3	Run 4	Run 5	Mean	SE
SG-BOED	Lower	5.537	5.536	5.473	5.523	5.518	<b>5.517</b>	<b>0.013</b>
SG-BOED	Upper	5.553	5.548	5.491	5.541	5.531	<b>5.533</b>	<b>0.012</b>
MINEBED	Lower	5.460	5.506	5.553	5.539	5.565	<b>5.524</b>	<b>0.021</b>
MINEBED	Upper	5.473	5.526	5.567	5.554	5.574	<b>5.540</b>	<b>0.022</b>

**Performance sensitivity to errors in the policy.** Finally, we investigate the effect of slight errors in the design policy network. To this end, we look at the performance achieved by partially trained design networks (there will be some errors or inaccuracies in networks that were not trained until convergence). Table 8 shows the performance of iDAD as a function of training time, demonstrating that small errors in the network only lead to small drops in performance.

In detail, our results show that with just 8% of the total training budget, this slightly inaccurate network still performs relatively well, achieving total mutual information of 7.1, compared to the fully trained network that reached 7.8. We also highlight that iDAD outperforms all baselines with as little as 1% of the total training budget (the best performing baseline achieves mutual information of 5.5, see Table 2).

## D.5 PK model

The drug concentration  $z$ , measured  $\xi$  hours after administering it, and the corresponding noisy observation  $y$  are given by

$$z(\xi; \theta) = \frac{D_V}{V} \frac{k_\alpha}{k_\alpha - k_e} [e^{-k_e \xi} - e^{-k_\alpha \xi}], \quad y(\xi; \theta) = z(\xi; \theta)(1 + \epsilon) + \eta \quad (105)$$

where  $\theta = (k_\alpha, k_e, V)$ ,  $D_V = 400$  is a constant,  $\epsilon \sim \mathcal{N}(0, 0.01)$  is multiplicative noise to account for heteroscedasticity and  $\eta \sim \mathcal{N}(0, 0.1)$  is an additive observation noise. Since both noise sources are Gaussian, the observation likelihood is also Gaussian i.e.

$$y(\xi; \theta) \sim \mathcal{N}(z(\xi; \theta), 0.01z(\xi; \theta)^2 + 0.1) \quad (106)$$

The prior for the parameters  $\theta$  that we used

$$\log \theta \sim \mathcal{N} \left( \begin{bmatrix} \log 1 \\ \log 0.1 \\ \log 20 \end{bmatrix}, \begin{bmatrix} 0.05 & 0 & 0 \\ 0 & 0.05 & 0 \\ 0 & 0 & 0.05 \end{bmatrix} \right) \quad (107)$$

### D.5.1 Training details

The architecture of the design network  $\pi_\phi$  used for Figure 3 and 4 and all its hyperparameters are in the following tables. For the encoder of the design-outcome pairs we used the following:



Layer	Description	iDAD, InfoNCE	iDAD, NWJ	Activation
Input	$\xi, y$	2	2	-
H1	Fully connected	64	64	ReLU
H2	Fully connected	512	512	ReLU
Output	Fully connected	32	32	-
Attention	8 heads	32	32	-

The outputs of the encoder,  $\{R(h_t)\}_{t=1}^T$ , are summed and the resulting vector (of dimension 32) is fed into an emitter network, for which we used the following:

Layer	Description	iDAD, InfoNCE	iDAD, NWJ	Activation
Input	$R(h_t)$	32	32	-
H1	Fully connected	256	256	ReLU
H2	Fully connected	32	32	ReLU
Output	Fully connected	1	1	Sigmoid

The architecture of the critic network  $U_\psi$  used in Figures 3 and 4 and all its hyperparameters are in the following tables. For the encoder of the design-outcome pairs we used the same architecture as for the design network, namely:

Layer	Description	iDAD, InfoNCE	iDAD, NWJ	Activation
Input	$\xi, y$	2	2	-
H1	Fully connected	64	64	ReLU
H2	Fully connected	512	512	ReLU
Output	Fully connected	32	32	-
Attention	8 heads	32	32	-

The resulting pooled representation,  $R(h_T)$  is fed into fully connected critic head layers with the following architecture:

Layer	Description	iDAD, InfoNCE	iDAD, NWJ	Activation
Input	$R(h_T)$	32	32	-
H1	Fully connected	512	512	ReLU
H2	Fully connected	256	256	ReLU
H3	Fully connected	512	512	ReLU
Output	Fully connected	32	32	-

Finally, for the latent variable encoder network we used:

Layer	Description	iDAD, InfoNCE	iDAD, NWJ	Activation
Input	$\theta$	3	3	-
H1	Fully connected	8	8	ReLU
H2	Fully connected	64	64	ReLU
H3	Fully connected	512	512	ReLU
Output	Fully connected	32	32	-

The optimisation was performed with Adam [26] with the following hyperparameters:

Parameter	iDAD, InfoNCE	iDAD, NWJ
Batch size	1024	1024
Number of contrastive/negative samples	1023	1023
Number of gradient steps	100000	100000
Initial learning rate (LR)	0.0001	0.0001
LR annealing factor	0.8	0.5
LR annealing frequency (if no improvement)	2000	2000

Table 9: Upper and lower bounds on the total information,  $\mathcal{I}_5(\pi)$ , for the pharmacokinetic experiment. Errors indicate  $\pm 1$  s.e. estimated over 4096 (126 for variational) histories and  $L = 5 \times 10^5$ .

Method	Lower bound	Upper bound	Deployment time
Random	$2.523 \pm 0.033$	$2.523 \pm 0.033$	N/A
Equal interval	$2.651 \pm 0.022$	$2.651 \pm 0.022$	N/A
MINEBED	$2.955 \pm 0.030$	$2.956 \pm 0.030$	N/A
SG-BOED	$2.985 \pm 0.027$	$2.985 \pm 0.027$	N/A
Variational	$2.683 \pm 0.093$	$2.683 \pm 0.093$	505.4 $\pm 1\%$
<b>IDAD (NWJ)</b>	<b><math>3.163 \pm 0.023</math></b>	<b><math>3.163 \pm 0.023</math></b>	0.007 $\pm 7\%$
<b>IDAD (InfoNCE)</b>	<b><math>3.200 \pm 0.024</math></b>	<b><math>3.200 \pm 0.024</math></b>	0.007 $\pm 8\%$
DAD	$3.234 \pm 0.023$	$3.234 \pm 0.023$	0.002 $\pm 7\%$

Table 10: Upper and lower bounds on the total information,  $\mathcal{I}_{10}(\pi)$ , for the pharmacokinetic experiment. Errors indicate  $\pm 1$  s.e. estimated over 4096 (126 for variational) histories and  $L = 5 \times 10^5$ .

Method	Lower bound	Upper bound	Deployment time
Random	$3.344 \pm 0.034$	$3.345 \pm 0.034$	N/A
Equal interval	$3.422 \pm 0.026$	$3.423 \pm 0.026$	N/A
MINEBED	$3.849 \pm 0.034$	$3.849 \pm 0.034$	N/A
SG-BOED	$3.824 \pm 0.034$	$3.824 \pm 0.034$	N/A
Variational	$3.624 \pm 0.099$	$3.624 \pm 0.099$	1055.2 $\pm 8\%$
<b>IDAD (NWJ)</b>	<b><math>4.034 \pm 0.025</math></b>	<b><math>4.034 \pm 0.025</math></b>	0.007 $\pm 6\%$
<b>IDAD (InfoNCE)</b>	<b><math>4.045 \pm 0.026</math></b>	<b><math>4.045 \pm 0.026</math></b>	0.007 $\pm 5\%$
DAD	$4.116 \pm 0.024$	$4.117 \pm 0.024$	0.007 $\pm 8\%$

## D.5.2 Hyperparameter selection

Hyperparameter selection was done in a way similar to the Location Finding experiment (see D.4.3). We tried encoding dimensions  $ED = 32, 64$  and selected the smaller size as there were no clear benefits to larger networks (relatively speaking, this is an easier model than the location finding). We used the same hidden dimension, i.e.  $HD = 512$ . In terms of learning rates, we tried 0.0001, 0.0005 and 0.001; we found 0.0001 to be appropriate, although NWJ bound was exhibiting high variance, so used a smaller learning rate annealing factor for that network (0.5 vs 0.8 for InfoNCE). We performed similar level of hyperparameter tuning for all trainable baselines as well (DAD, MINEBED and SG-BOED).

## D.5.3 Further results

Table 9 reports the results shown in Figure 3c), along with the corresponding upper bounds and deployment times, while Table 10 reports the results for  $T = 10$ .

**Training stability.** To assess the robustness of the results and the stability of the training process, we trained 5 additional iDAD networks with each of the two bounds, using different seeds but the same hyperparameters we used to produce the results of the pharmacokinetic experiment (Figure 3c) and corresponding Table 9). We report upper and lower bounds on the mutual information along with their mean and standard error in the table below.

Method	Bound	Run 1	Run 2	Run 3	Run 4	Run 5	Mean	SE
iDAD, InfoNCE	Lower	3.209	3.165	3.198	3.221	3.128	<b>3.185</b>	<b>0.019</b>
iDAD, InfoNCE	Upper	3.210	3.166	3.201	3.223	3.130	<b>3.186</b>	<b>0.019</b>
iDAD, NWJ	Lower	3.034	3.049	2.608	3.149	3.082	<b>3.034</b>	<b>0.107</b>
iDAD, NWJ	Upper	3.034	3.049	2.609	3.150	3.083	<b>3.034</b>	<b>0.107</b>

We repeat the same procedure for the static baselines. The results reported in the table below demonstrate the training stability of these baselines as well.

Method	Bound	Run 1	Run 2	Run 3	Run 4	Run 5	Mean	SE
SG-BOED	Lower	2.932	2.452	2.448	2.991	2.962	<b>2.757</b>	<b>0.140</b>
SG-BOED	Upper	2.932	2.453	2.449	2.992	2.962	<b>2.757</b>	<b>0.140</b>
MINEBED	Lower	2.912	2.213	3.014	2.092	2.941	<b>2.634</b>	<b>0.221</b>
MINEBED	Upper	2.914	2.213	3.015	2.092	2.942	<b>2.635</b>	<b>0.222</b>

## D.6 SIR Model

Generally speaking, the SIR model advocates that, within a fixed population of size  $N$ , susceptible individuals  $S(\tau)$ , where  $\tau$  is time, can become infected and move to an infected state  $I(\tau)$ . The infected individuals can then recover from the disease and move to the recovered state  $R(\tau)$ . The dynamics of these events are governed by the infection rate  $\beta$  and recovery rate  $\gamma$ , which define the particular disease in question. In the context of BOED, the aim is generally to estimate these two model parameters by observing state populations at particular measurement times  $\tau$ , which are the experimental design variables. The SIR model has been studied extensively in the context of BOED, e.g. in [12, 27, 29, 30].

Stochastic versions of the SIR model are usually formulated via continuous-time Markov chains (CTMC), which can be simulated from via the Gillespie algorithm [2], yielding discrete state populations. However, iDAD requires us to differentiate through the sampling path of the state populations to the experimental designs, which is impossible if the simulated data is discrete as gradients are undefined. Thus, we here implement an alternative formulation of the stochastic SIR model that is based on stochastic differential equations (SDEs), as studied in [29], which yields continuous state populations that can be differentiated.

Following [29], let us first define a state population vector  $\mathbf{X}(\tau) = (S(\tau), I(\tau))^\top$ , where we can safely ignore the population of recovered  $R(\tau)$  for modelling purposes because we assume that the total population stays fixed. The system of Itô SDEs that defines the stochastic SIR model is given by

$$d\mathbf{X}(\tau) = \mathbf{f}(\mathbf{X}(\tau))d\tau + \mathbf{G}(\mathbf{X}(\tau))d\mathbf{W}(\tau), \quad (108)$$

where  $\mathbf{f}$  is a drift vector,  $\mathbf{G}$  is a diffusion matrix and  $\mathbf{W}(\tau)$  is a vector of independent Wiener processes. [29] showed that the drift vector and diffusion matrix are given by

$$\mathbf{f}(\mathbf{X}(\tau)) = \begin{pmatrix} -\beta \frac{S(\tau)I(\tau)}{N} \\ \beta \frac{S(\tau)I(\tau)}{N} - \gamma I(\tau) \end{pmatrix} \quad \text{and} \quad \mathbf{G}(\mathbf{X}(\tau)) = \begin{pmatrix} -\sqrt{\beta \frac{S(\tau)I(\tau)}{N}} & 0 \\ \sqrt{\beta \frac{S(\tau)I(\tau)}{N}} & -\sqrt{\gamma I(\tau)} \end{pmatrix}. \quad (109)$$

Given the system of Itô SDEs in (108), as well as the above drift vector and diffusion matrix, we can then simulate state populations  $\mathbf{X}(\tau)$  by solving the SDE using finite-difference methods, such as e.g. the Euler-Maruyama method. See [29] for more information on the SDE-based SIR model, including derivations of the drift vector and diffusion matrix.

Importantly, we note that [29] further used the solutions of (108) as an input to a Poisson observation model, which increases the noise in simulated data. We here opt to simply use the solutions of (108) as data and do not consider an additional Poisson observational model.

### D.6.1 Training details

As previously mentioned, the design variable for this model is the measurement time  $\tau \in [0, 100]$ . When solving the SDE with the Euler-Maruyama method, we discretize the time domain with a resolution of  $\Delta\tau = 10^{-2}$ . We here only use the number of infected  $I(\tau)$  as the observed data, as others might be difficult to measure in reality. The total population is fixed at  $N = 500$  and the initial conditions are  $\mathbf{X}(\tau = 0) = (0, 2)^\top$ . The model parameters  $\beta$  and  $\gamma$  have log-normal priors, i.e.  $p(\beta) = \text{Lognorm}(0.50, 0.50^2)$  and  $p(\gamma) = \text{Lognorm}(0.10, 0.50^2)$ . Importantly, because solving SDEs is expensive, we pre-simulate our data on a time grid, store it in memory and then access the relevant data during training.

We present the network architectures and hyper-parameters corresponding to the  $T = 5$  iDAD results shown in Table 5 of the main text. For the encoder of the design-outcome pairs we used:

Layer	Description	iDAD, InfoNCE	iDAD, NWJ	Activation
Input	$\xi, y$	2	2	-
H1	Fully connected	8	8	ReLU
H2	Fully connected	64	64	ReLU
H3	Fully connected	512	512	ReLU
Output	Fully connected	32	32	-

The resulting representations,  $\{R(h_t)\}_{t=1}^{T-1}$ , are stacked into a matrix (as new design–outcome pairs are obtained) and fed into an emitter network, which contains an LSTM cell with two hidden layers. We only keep the last hidden state of the LSTM’s output and pass it through a final FC layer:

Layer	Description	iDAD, InfoNCE	iDAD, NWJ	Activation
Input	$\{R(h_t)\}_{t=1}^{T-1}$	$32 \times t$	$32 \times t$	-
H1 & H2	LSTM	32	32	-
H3	Fully connected	16	16	ReLU
Output	Fully connected	1	1	-

The architecture of the critic network  $U_\psi$  used in Table 5 and all its hyper-parameters are in the tables that follow. First, the encoder network of the latent variables is:

Layer	Description	iDAD, InfoNCE	iDAD, NWJ	Activation
Input	$\theta$	2	2	-
H1	Fully connected	8	8	ReLU
H2	Fully connected	64	64	ReLU
H3	Fully connected	512	512	ReLU
Output	Fully connected	32	32	-

For the design–outcome pairs encoder we use the same architecture as in the design network, namely:

Layer	Description	iDAD, InfoNCE	iDAD, NWJ	Activation
Input	$\xi, y$	2	2	-
H1	Fully connected	8	8	ReLU
H2	Fully connected	64	64	ReLU
H3	Fully connected	512	512	ReLU
Output	Fully connected	32	32	-

The outputs of the encoder,  $\{R(h_t)\}_t$ , are stacked and fed into an LSTM cell with two hidden layers. We only keep the last hidden state of the LSTM’s output and pass it through a FC layer:

Layer	Description	iDAD, InfoNCE	iDAD, NWJ	Activation
Input	$\{R(h_t)\}_{t=1}^{T-1}$	$32 \times t$	$32 \times t$	-
H1 & H2	LSTM	32	32	-
H3	Fully connected	16	16	ReLU
Output	Fully connected	32	32	-

The optimization was performed with Adam [26] with learning rate annealing with the following hyper-parameters:

Parameter	iDAD InfoNCE	iDAD, NWJ
Batch size	512	512
Number of contrastive/negative samples	511	511
Number of gradient steps	100000	100000
Initial learning rate (LR)	0.0005	0.0005
LR annealing factor	0.96	0.96
LR annealing frequency	1000	1000

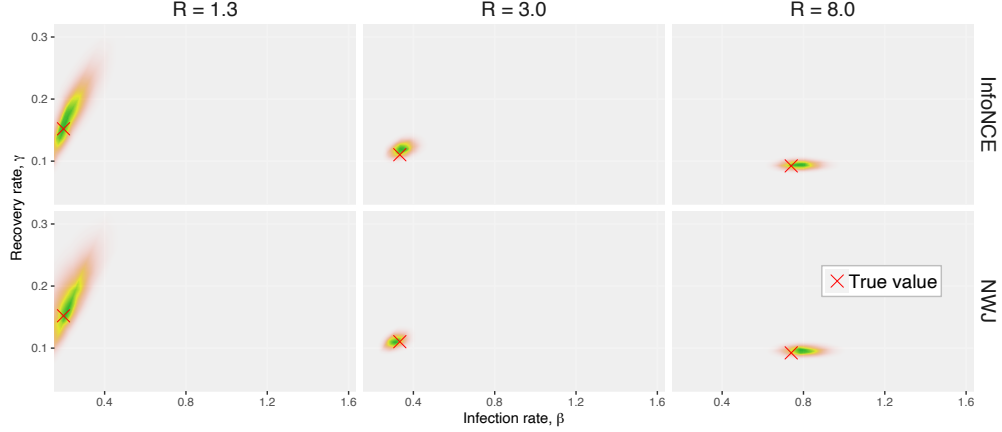


Figure 8: Approximate posteriors for the SIR model.

### D.6.2 Further results

**Different number of experiments  $T$ .** In Table 11 we show lower bound estimates when applying iDAD with the InfoNCE lower bound to the SDE-based SIR model for different number of measurements  $T$ . The design network and critic architectures are the same as for  $T = 5$ . Table 11 shows that more measurements yield higher expected information gains, as one might intuitively expect. Furthermore, the increase in expected information gain saturates with increasing  $T$ , which is why we presented the results for  $T = 5$  in the main text. The biggest increase, however, occurs from  $T = 1$  to  $T = 2$ . This is intuitive, because the SIR model has two model parameters that we wish to estimate but we only gather one data point with one measurement. Hence, in order to accurately estimate both of these parameters, we would need at least 2 measurements, which is reflected in Table 11. We note that all of these numbers, with the exception of  $T = 1$ , are larger than those found by [29]. This increase in expected information gain may be explained by the fact that [29] use an additional Poisson observation model, which means that the resulting data are inherently noisier and less informative.

Table 11: InfoNCE lower bound estimates ( $\pm$  s.e.) when applying iDAD to the SDE-based SIR model for different number of measurements  $T$ .

$T$	iDAD, InfoNCE	iDAD, NWJ
1	$1.396 \pm 0.018$	$1.417 \pm 0.001$
2	$2.714 \pm 0.019$	$2.699 \pm 0.001$
3	$3.554 \pm 0.021$	$3.515 \pm 0.001$
4	$3.600 \pm 0.018$	$3.749 \pm 0.001$
5	$3.915 \pm 0.020$	$3.869 \pm 0.001$
7	$4.027 \pm 0.019$	$3.911 \pm 0.001$
10	$4.100 \pm 0.020$	$4.019 \pm 0.001$

**Training stability.** To assess the robustness of the results and the stability of the training process, we trained 5 additional iDAD networks with each of the two bounds, using different seeds but the same hyperparameters we used to produce the results of Table 5 in the main text. We report upper and lower bounds on the mutual information along with their mean and standard error in the table below.

Method	Bound	Run 1	Run 2	Run 3	Run 4	Run 5	Mean	SE
iDAD, InfoNCE	Lower	3.900	3.919	3.919	3.901	3.887	<b>3.906</b>	<b>0.007</b>
iDAD, NWJ	Lower	3.872	3.838	3.854	3.883	3.848	<b>3.859</b>	<b>0.009</b>

We repeat the same procedure for the static baselines. The results reported in the table below demonstrate the training stability of these baselines as well.

Method	Bound	Run 1	Run 2	Run 3	Run 4	Run 5	Mean	SE
SG-BOED	Lower	3.713	3.765	3.767	3.764	3.739	<b>3.749</b>	<b>0.012</b>
MINEBED	Lower	3.373	3.438	3.376	3.379	3.420	<b>3.397</b>	<b>0.015</b>