

---

# Hallucinated Adversarial Control for Conservative Offline Policy Evaluation

---

Jonas Rothfuss\*<sup>1</sup>

Bhavya Sukhija\*<sup>1</sup>

Tobias Birchler\*<sup>1</sup>

Parnian Kassarai<sup>1</sup>

Andreas Krause<sup>1</sup>

<sup>1</sup>ETH Zurich, Switzerland

## Abstract

We study the problem of *conservative off-policy evaluation (COPE)* where given an offline dataset of environment interactions, collected by other agents, we seek to obtain a (tight) lower bound on a policy’s performance. This is crucial when deciding whether a given policy satisfies certain minimal performance/safety criteria before it can be deployed in the real world. To this end, we introduce HAMBO, which builds on an uncertainty-aware learned model of the transition dynamics. To form a conservative estimate of the policy’s performance, HAMBO hallucinates worst-case trajectories that the policy may take, within the margin of the models’ epistemic confidence regions. We prove that the resulting COPE estimates are valid lower bounds, and, under regularity conditions, show their convergence to the true expected return. Finally, we discuss scalable variants of our approach based on Bayesian Neural Networks and empirically demonstrate that they yield reliable and tight lower bounds in various continuous control environments.

## 1 INTRODUCTION

Reinforcement learning methods require many interactions with their environment to successfully learn and evaluate policies. Therefore, they are rarely applied in challenging real-world applications such as medicine [Murphy et al., 2001], education [Mandel et al., 2014] or autonomous driving [Kiran et al., 2021], where a policy can only be deployed in the environment if it exceeds a pre-specified performance threshold or fulfills certain safety criteria. This leaves us with a challenging problem: How do we know whether a policy fulfills the necessary criteria so that it can

safely interact with the environment, without testing it on the environment, and in the process, compromising safety?

Off-policy evaluation (OPE) aims to solve this problem by estimating the performance of an evaluation policy, using only offline data that was previously collected by other agents [e.g. Precup et al., 2001, Dudík et al., 2011]. In practice, offline datasets are often recorded interactions of a human expert with the environment. Since the evaluation policy typically induces a different action-state distribution than offline data, OPE methods often have to make predictions under strong distribution shifts. As a result, most existing OPE estimators suffer from high variance and are prone to overestimating the performance of the policy [Thomas et al., 2015]. In safety-critical applications, we can not risk and deploy a policy that is potentially much worse than what the OPE estimate suggests. Therefore, we aim for *conservative off-policy evaluation (COPE)* which seeks a (tight) lower bound on the evaluation policy’s expected return that holds with high probability. Once deployed, the policy may end up exploring areas that were not included in the offline data. Thus, reliably bounding the worst-case performance can be quite challenging.

We develop a novel *model-based* COPE approach that hinges upon two key ideas: *epistemic uncertainty* and *pessimism*. In particular, our approach, *Hallucinated Adversarial Model-Based Off-policy evaluation (HAMBO)* (HAMBO), builds on a learned statistical model of the transition dynamics that is able to quantify epistemic uncertainty. To obtain a valid lower bound on the policy performance, HAMBO hallucinates adversarial/worst-case trajectories the agent may take within the epistemic confidence sets of the model.

We prove that HAMBO reliably yields a high-probability bound on the true expected return of the policy, even when the offline data does not cover the areas explored by the evaluation policy (Proposition 3.2). Under regularity conditions, we further show that our conservative estimate *converges* from below to the true expected return (Theo-

\*Equal contribution.

rem 3.8). To the best of our knowledge, HAMBO is the first provably consistent and conservative approach for OPE in continuous action-state spaces. We then propose scalable Bayesian neural network (BNN) variants of HAMBO and empirically evaluate them on various continuous control tasks. Importantly, we demonstrate that, *even when the regularity conditions are not met*, HAMBO reliably provides tight lower bounds on the true expected return.

## 2 PROBLEM SETTING

We consider a finite horizon Markov decision process (MDP)  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, p_0, p, r, T)$  with continuous state and action spaces  $\mathcal{S} \subseteq \mathbb{R}^{d_s}$  and  $\mathcal{A} \subseteq \mathbb{R}^{d_a}$ , initial state distribution  $p_0(s_0)$ , reward function  $r(\mathbf{a}_t, \mathbf{s}_t)$  and horizon  $T \in \mathbb{N}$ . In particular, we consider stochastic transition dynamics that are governed by  $\mathbf{s}_{t+1} = \mathbf{f}(\mathbf{s}_t, \mathbf{a}_t) + \boldsymbol{\epsilon}_t$  where  $\mathbf{f} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$  is unknown and  $\boldsymbol{\epsilon}_t \in \mathbb{R}^{d_s}$  is independent, additive transition noise with distribution  $p_\epsilon(\boldsymbol{\epsilon}_t | \mathbf{s}_t, \mathbf{a}_t)$ . Hence, the transition distribution  $p$  follows as  $p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) = p_\epsilon(\mathbf{s}_{t+1} - \mathbf{f}(\mathbf{s}_t, \mathbf{a}_t) | \mathbf{s}_t, \mathbf{a}_t)$ . For simplicity, we assume that the reward function is known. However, all results can straightforwardly be extended to unknown rewards.

The agent interacts with the environment according to a policy  $\pi(\mathbf{a}_t | \mathbf{s}_t)$ , which is a distribution over actions, conditioned on the current state  $\mathbf{s}_t$ . The performance of a policy is typically measured by its expected return  $J(\pi) := J_p(\pi) := \mathbb{E}_{s_0 \sim p_0}[V_{p,0}^\pi(s_0)]$  where  $V_t^\pi(\mathbf{s}) := V_{p,t}^\pi(\mathbf{s}) := \mathbb{E}_{p,\pi}[G_t | S_t = \mathbf{s}]$  is the value function and  $G_t := \sum_{t'=t+1}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'})$  is the return. For simplicity, we omit a discount factor in the return computation. However, all results presented can be straightforwardly adapted to discounted rewards. Furthermore, we denote the *occupancy measure* of policy  $\pi$  as

$$\rho^\pi(\mathbf{s}, \mathbf{a}) := \frac{1}{T} \sum_{t=0}^{T-1} p(\mathbf{s}_t = \mathbf{s}, \mathbf{a}_t = \mathbf{a} | \pi, \mathcal{M}),$$

that is, the probability density function of being in state  $\mathbf{s}$  and performing action  $\mathbf{a}$  at any point of time  $t = 0, \dots, T-1$ .

We study the problem of offline policy evaluation where the task is to evaluate the performance, i.e. estimate the expected return  $J(\pi_e)$ , of a given evaluation policy  $\pi_e$  while only using an offline dataset  $\mathcal{D}_b = \{(\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i)\}_{i=1}^n$  of observed transitions. The key challenge in OPE is the distribution shift between the (unknown) behavior policy  $\pi_b$  which generated the dataset  $\mathcal{D}_b$  and the policy  $\pi_e$  which we would like to evaluate. If  $\pi_b$  differs from  $\pi_e$ , their state occupancy measures  $\rho^{\pi_b}$  and  $\rho^{\pi_e}$  can look significantly different. As a result, the dataset  $\mathcal{D}_b$  which is generated based on  $\rho^{\pi_b}$  may contain many samples in regions of the state-action space which  $\pi_e$  is unlikely to visit and limited data in regions that are relevant for accurately evaluating  $\pi_e$ . In some cases, the support of  $\rho^{\pi_b}$  might not even contain the support of

$\rho^{\pi_e}$ , i.e.,  $\exists(\mathbf{s}, \mathbf{a}) \in \mathcal{S} \times \mathcal{A} : \rho^{\pi_e}(\mathbf{s}, \mathbf{a}) > 0 \wedge \rho^{\pi_b}(\mathbf{s}, \mathbf{a}) = 0$ . Since OPE methods have to make predictions under such strong distribution shifts their estimates suffer from high variance and are prone to overestimate the performance of the policy.

OPE is particularly relevant in applications where we need to ensure a certain level of performance before a policy can be deployed online. Hence, it is often important to reliably determine whether or not the policy  $\pi_e$  meets its minimum performance requirements. We formalize this problem as *conservative offline policy evaluation* (see Definition 2.1) where we want to ideally find a tight lower bound on the expected return that holds with high-probability:

**Definition 2.1** (Conservative Offline Policy Evaluation). *Let  $\mathcal{M}$  be an MDP and  $\mathcal{D}_b \in (\mathcal{S} \times \mathcal{A} \times \mathbb{R} \times \mathcal{S})^n$  a dataset of transitions, collected with a behavior policy  $\pi_b$  on  $\mathcal{M}$ . Then the task of conservative OPE is: Given the offline dataset  $\mathcal{D}_b$ , a policy  $\pi_e$  to evaluate and a confidence level  $\delta \in (0, 1)$ , find the largest possible lower-bound  $b \in \mathbb{R}$ , which satisfies  $b \leq J(\pi_e)$  with probability at least  $1 - \delta$ .*

In some applications [e.g., Brunke et al., 2022], safety criteria are not directly encoded in the reward and instead, are expressed as additional constraints in the form of  $\mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \rho^{\pi_e}} [c_i(\mathbf{s}, \mathbf{a})] \geq 0$ . To determine with high confidence whether  $\pi_e$  meets these constraints, we can apply COPE to each  $c_i$  individually.

## 3 COPE VIA ADVERSARIAL TRANSITION MODELS

We take a model-based approach to COPE, and use a statistical model to estimate which transition functions  $\mathbf{h} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$  from a hypothesis space  $\mathcal{H}$  are plausible given the offline data  $\mathcal{D}_b$  of size  $n$ . Then, we employ this statistical model of the transition dynamics to estimate the policy value  $J(\pi_e)$ . For this estimate, we want to guarantee with high probability that it does not exceed the true policy value. To ensure this, we need to be able to reliably quantify the *epistemic uncertainty* of our model estimates.

Uncertainty quantification can be done with either a frequentist approach that produces mean and confidence estimate  $\boldsymbol{\mu}_n(\mathbf{s}, \mathbf{a})$  and  $\boldsymbol{\sigma}_n(\mathbf{s}, \mathbf{a})$  or with a Bayesian model that maintains a posterior distribution  $p(\mathbf{h} | \mathcal{D}_b)$  over dynamics models in  $\mathcal{H}$ . In the Bayesian case, we denote  $\boldsymbol{\mu}_n(\mathbf{s}, \mathbf{a}) := \mathbb{E}_{\mathbf{h} \sim p(\mathbf{h} | \mathcal{D}_b)}[\mathbf{h}(\mathbf{s}, \mathbf{a})]$  as the posterior mean and  $\boldsymbol{\sigma}_n^2(\mathbf{s}, \mathbf{a}) := \text{diag}(\mathbb{E}_{\mathbf{h}, \mathbf{h}' \sim p(\mathbf{h} | \mathcal{D}_b)}[\mathbf{h}(\mathbf{s}, \mathbf{a}) \mathbf{h}'(\mathbf{s}, \mathbf{a})^\top])$  as the posterior variance. In either case, we require that our statistical model of  $\mathbf{h}$  is calibrated:

**Assumption 3.1** (Calibrated model). *A statistical model  $(\boldsymbol{\mu}_n, \boldsymbol{\sigma}_n, \beta_n)$ , with  $\beta_n(\delta) \in \mathbb{R}^+$  as a scalar function that depends on the confidence level  $\delta \in (0, 1]$ , is calibrated*

with respect to  $\mathbf{f}$  if, with probability at least  $1 - \delta$ , for all  $(\mathbf{s}, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}$  and  $j = 1, \dots, d_s$

$$|\mu_{n,j}(\mathbf{s}, \mathbf{a}) - f_j(\mathbf{s}, \mathbf{a})| \leq \beta_n(\delta)\sigma_{n,j}(\mathbf{s}, \mathbf{a}),$$

where  $\mu_{n,j}$  and  $\sigma_{n,j}$  denote the  $j$ -th element in the vector-valued functions  $\boldsymbol{\mu}_n$  and  $\boldsymbol{\sigma}_n$ , respectively.

Popular statistical models for transition dynamics that capture epistemic uncertainty are *Gaussian Processes (GPs)* [Rasmussen and Williams, 2005], *Probabilistic Neural Network Ensembles* [Lakshminarayanan et al., 2017] and *Bayesian Neural Networks* [Blundell et al., 2015]. In later sections, we will attend to these specific choices of model in more detail and discuss when they are calibrated.

### 3.1 THE HAMBO FRAMEWORK

If our model is calibrated, we can, with high probability, use the confidence region

$$[\boldsymbol{\mu}_n(\mathbf{s}, \mathbf{a}) - \beta_n(\delta)\boldsymbol{\sigma}_n(\mathbf{s}, \mathbf{a}), \boldsymbol{\mu}_n(\mathbf{s}, \mathbf{a}) + \beta_n(\delta)\boldsymbol{\sigma}_n(\mathbf{s}, \mathbf{a})]$$

which is a  $d_s$ -dimensional hypercube, as a proxy for the true dynamics  $\mathbf{f}(\mathbf{s}, \mathbf{a})$ . We then pessimistically select transitions within this region, to guarantee a high probability lower bound on the policy value  $J(\pi_e)$ . We do so, by introducing an adversary  $\boldsymbol{\eta} : \mathcal{S} \times \mathcal{A} \rightarrow [-1, 1]^{d_s}$  that, for every  $(\mathbf{s}, \mathbf{a}) \in \mathcal{S} \times \mathcal{A}$  picks a transition from the confidence region, thereby inducing the following hallucinated transition distribution:

$$\tilde{p}_\eta(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t) := p_\epsilon(\mathbf{s}_{t+1} - \boldsymbol{\mu}_n(\mathbf{s}_t, \mathbf{a}_t) - \beta_n\boldsymbol{\eta}(\mathbf{s}_t, \mathbf{a}_t)\boldsymbol{\sigma}_n(\mathbf{s}_t, \mathbf{a}_t)). \quad (1)$$

This allows us to obtain a conservative value estimate for  $\pi_e$

$$\tilde{J}(\pi_e) := \min_{\boldsymbol{\eta}} J_{\tilde{p}_\eta}(\pi_e). \quad (2)$$

This equation summarizes our approach *hallucinated adversarial model-based off-policy evaluation (HAMBO)* and Algorithm 1 presents the pseudo-code. Here, the expected reward  $J_{\tilde{p}_\eta}(\pi_e)$  of  $\pi_e$  under the hallucinated transition model  $\tilde{p}_\eta$  can, e.g., be estimated via Monte Carlo estimation (i.e., generating trajectory rollouts and averaging the respecting returns). To find the adversary  $\boldsymbol{\eta}(\mathbf{s}, \mathbf{a})$  which minimizes (2), we can view  $\boldsymbol{\eta}(\mathbf{s}, \mathbf{a})$  as policy that aims to maximize  $-J_{\tilde{p}_\eta}(\pi_e)$  and solve the corresponding optimal control problem. Importantly, with high probability,  $\tilde{J}(\pi_e)$  is a lower bound on the true policy value  $J(\pi_e)$ :

**Proposition 3.2 (Valid lower bound).** *Given a calibrated model  $(\boldsymbol{\mu}_n, \boldsymbol{\sigma}_n, \beta_n(\delta))$ , the HAMBO estimates satisfy  $\tilde{J}(\pi_e) \leq J(\pi_e)$ , with probability greater than  $1 - \delta$ .*

While Proposition 3.2 shows that our estimate  $\tilde{J}(\pi_e)$  fulfills the requirements of COPE,  $\tilde{J}(\pi_e)$  could potentially be very loose. However, we can further establish a worst-case lower bound on  $\tilde{J}(\pi_e)$ , if  $\mathbf{f}$ ,  $r$ ,  $\boldsymbol{\sigma}_n$  and  $\pi_e$  are continuous. Formally, we make the following Lipschitz continuity assumption:

**Assumption 3.3. (Lipschitz continuity)**  $\mathbf{f}$  is  $L_f$ -Lipschitz,  $r$  is  $L_r$ -Lipschitz,  $\boldsymbol{\sigma}$  is  $L_\sigma$ -Lipschitz and  $\pi_e$  is  $L_\pi$ -Lipschitz w.r.t. the Wasserstein-1 distance, i.e., for all  $\mathbf{s}, \mathbf{s}' \in \mathcal{S}$

$$\mathcal{W}_1(\pi(\mathbf{a}|\mathbf{s}), \pi(\mathbf{a}|\mathbf{s}')) \leq L_\pi \|\mathbf{s} - \mathbf{s}'\|_2. \quad (3)$$

Here, the continuity assumption on  $\pi$  is expressed in terms the Wasserstein-1 distance and implies that a small change in the state space only induces a proportionally small change in the conditional action distribution of the policy. For instance, this is the case for policies that can be reparametrized with a Lipschitz function which is very common in practice:

**Example 3.4.** *Any policy  $\pi(\mathbf{a}|\mathbf{s})$  that can be reparametrized as  $\mathbf{g}(\mathbf{s}, \zeta)$ , where  $\zeta \sim p(\zeta)$  and  $\mathbf{g}$  is  $L_g$ -Lipschitz, is also  $L_g$ -Lipschitz w.r.t. the  $\mathcal{W}_1$ -distance.*

Such Lipschitz assumptions are common in model-based OPE [e.g. Fonteneau et al., 2009, Paduraru, 2013] and RL more broadly [e.g. Berkenkamp et al., 2017, Curi et al., 2020], and, e.g, hold in many real-world control problems. With these regularity assumptions, we bound how far away the HAMBO estimate  $\tilde{J}(\pi_e)$  is from the true policy value:

**Theorem 3.5.** *Under Assumption 3.1 and 3.3 we have, with probability at least  $1 - \delta$ , that*

$$J(\pi_e) - \tilde{J}(\pi_e) \leq C_n \mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \rho^{\pi_e}} [\|\boldsymbol{\sigma}_n(\mathbf{s}, \mathbf{a})\|_2]$$

where

$$C_n := \bar{L}_r \left(1 + \sqrt{d_s}\right) \beta_n T^2 \left(1 + \bar{L}_f + (1 + \sqrt{d_s})\beta_n(\delta)\bar{L}_\sigma\right)^{T-1}$$

with  $\bar{L}_r := L_r(1 + L_\pi)$  and  $\bar{L}_f, \bar{L}_\sigma$  defined analogously.

This theorem shows how by tuning the confidence level  $\delta$ , we can trade-off accuracy with reliability. In particular, choosing a small  $\delta$  will ensure that the upper-bound on  $J(\pi_e)$  holds. However, it also increases  $\beta_n(\delta)$  and loosens the bound, indicating the  $\tilde{J}(\pi_e)$  estimate will be less accurate. Tightness of the HAMBO lower bound  $\tilde{J}(\pi_e)$  depends on the following key factors: Lipschitz-regularity, episode horizon  $T$ , and epistemic uncertainty. Mainly, smaller Lipschitz constants and shorter episode lengths improve the bound. Moreover, the smaller the expected epistemic standard deviation  $\boldsymbol{\sigma}_n(\mathbf{s}_t, \mathbf{a}_t)$  under the state occupancy measure of  $\pi_e$ , the tighter the bound. While the first two factors are generally dictated by the problem instance, the epistemic uncertainty can be reduced by using more offline data (in the relevant areas of the state-action space). If we can show that the epistemic uncertainty shrinks sufficiently fast with the number of offline data points  $n$  (i.e., faster than  $\mathcal{O}(\beta_n^T)$ ), then we can prove that  $\tilde{J}(\pi_e)$  converges to the true policy value as  $n \rightarrow \infty$ . In the following, we discuss corresponding sufficient convergence conditions for GP models.

### 3.2 HAMBO WITH SMOOTH GP FUNCTIONS

In this section, we discuss the application of GPs for constructing calibrated confidence regions to be used for HAMBO. For the transition dynamics, we consider vector-valued functions  $\mathbf{f}(s, \mathbf{a}) \mapsto (f_1(s, \mathbf{a}), \dots, f_{d_s}(s, \mathbf{a}))$  such that the scalar-valued functions  $f_j \in \mathcal{H}_k$  reside in a Reproducing Kernel Hilbert Space (RKHS)  $\mathcal{H}_k$  with kernel function  $k(\cdot, \cdot)$  and have bounded RKHS norm, i.e.  $\|f_j\|_k \leq B$ . We denote this space by  $\mathbf{f} \in \mathcal{H}_{k,B}^{d_s} = \{[f_1, \dots, f_{d_s}] : \|f_j\|_k \leq B, j = 1, \dots, d_s\}$ . We assume that transition noise  $\epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2 \mathbf{I})$  is normally distributed with variance  $\sigma_\epsilon^2$ .

By fitting a zero-mean Gaussian Process  $\mathcal{GP}(0, k)$  on each dimension  $j = 1, \dots, d_s$  of the next state  $s_{t+1}$ , we can use the posterior means and variances to construct calibrated confidence sets. For brevity, we denote  $\mathbf{x} := (s, \mathbf{a})$ , so that

$$\begin{aligned} \mu_{n,j}(\mathbf{x}) &= \mathbf{k}_n^\top(\mathbf{x})(\mathbf{K}_n + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{y}_{n,j} \\ \sigma_{n,j}^2(\mathbf{x}) &= k(\mathbf{x}, \mathbf{x}) - \mathbf{k}_n^\top(\mathbf{x})(\mathbf{K}_n + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{k}_n(\mathbf{x}) \end{aligned} \quad (4)$$

where  $\mathbf{y}_{n,j} = [s'_{i,j}]_{i \leq n}^\top$  is the vector the  $j$ -th element of the observed next states  $s'_i$ ,  $\mathbf{k}_n(\mathbf{x}) = [k(\mathbf{x}, \mathbf{x}_i)]_{i \leq n}^\top$ , and  $\mathbf{K}_n = [k(\mathbf{x}_i, \mathbf{x}_l)]_{i,l \leq n}$  is the kernel matrix. By concatenating the element-wise posterior mean and standard deviation, we obtain  $\boldsymbol{\mu}_n(\mathbf{x}) = [\mu_{n,j}(\mathbf{x})]_{j \leq d_s}^\top$  and  $\boldsymbol{\sigma}_n(\mathbf{x}) = [\sigma_{n,j}(\mathbf{x})]_{j \leq d_s}^\top$ . Using this, we can construct calibrated confidence intervals that fulfill Assumption 3.1:

**Lemma 3.6** (Calibrated GP confidence sets). *Let  $\mathbf{f} \in \mathcal{H}_{k,B}^{d_s}$ . Suppose  $\boldsymbol{\mu}_n$  and  $\boldsymbol{\sigma}_n$  are the posterior mean and variance of a GP with kernel  $k$ , fitted to  $n$  noisy evaluations of  $\mathbf{f}$ . There exists  $\beta_n(\delta)$ , for which the tuple  $(\boldsymbol{\mu}_n, \boldsymbol{\sigma}_n, \beta_n(\delta))$  satisfies Assumption 3.1 w.r.t. function  $\mathbf{f}$ .*

In Appendix B.2 we prove this lemma using results of Chowdhury and Gopalan [2017] and give the exact expression for a  $\beta_n(\delta)$  that satisfies it. Generally,  $\beta_n(\delta)$  depends on the maximum information capacity  $\gamma_n$  of the kernel (see Appendix B.2 for definition and details). In the GP setting, we can also show Lipschitz continuity of  $\mathbf{f}$  and  $\boldsymbol{\sigma}$ , if the kernel function  $k$  is sufficiently regular:

**Lemma 3.7.** *If the kernel metric  $d_k(\mathbf{x}, \mathbf{x}') := (k(\mathbf{x}, \mathbf{x}) + k(\mathbf{x}', \mathbf{x}') - 2k(\mathbf{x}, \mathbf{x}'))^{\frac{1}{2}}$  is  $L_k$ -Lipschitz, then every  $\mathbf{f} \in \mathcal{H}_{k,B}^{d_s}$  is Lipschitz with  $L_f = \sqrt{d_s} B L_k$  and the posterior standard deviation  $\boldsymbol{\sigma}$  is Lipschitz with  $L_\sigma = \sqrt{d_s} L_k$ .*

For common kernels, the kernel metric is Lipschitz continuous, and thus Lemma 3.7 applies. For instance, for the linear kernel we have  $L_k = 1$ , for the RBF kernel we have  $L_k = 1/\ell$  and for the Matern- $\nu$  kernel we have  $L_k = \sqrt{\nu/(\nu-1)}/\ell$ , where  $\ell$  is the lengthscale and  $\nu$  the smoothness parameter of the Matern kernel.

We can conclude that conditions of Proposition 3.2 and Theorem 3.5 are met when a GP is used for learning the tran-

sition dynamics from offline data. Hence, when the reward and the policy are Lipschitz, the HAMBO estimate satisfies

$$J(\pi_\epsilon) - C_n \mathbb{E}_{\rho^{\pi_\epsilon}} [\|\boldsymbol{\sigma}_n(s, \mathbf{a})\|_2] \leq \tilde{J}(\pi_\epsilon) \leq J(\pi_\epsilon)$$

with high probability. We can show that given a dataset of i.i.d. trajectories, the difference term shrinks with  $n$  sufficiently fast:

**Theorem 3.8** (Consistency of HAMBO). *Let  $r$  be  $L_r$ -Lipschitz,  $\pi$  be  $L_\pi$ -Lipschitz w.r.t. the  $\mathcal{W}_1$ -distance and  $\mathbf{f} \in \mathcal{H}_{k,B}^{d_s}$  where  $k$  is a kernel with a  $L_k$ -Lipschitz kernel metric with a maximum information capacity  $\gamma_n$  which is  $\mathcal{O}(\text{polylog}(n))$ . Suppose both  $\rho^{\pi_\epsilon}$  and  $\rho^{\pi_b}$  have a compact support and  $\text{supp}(\rho^{\pi_\epsilon}) \subseteq \text{supp}(\rho^{\pi_b})$  and  $\mathcal{D}_b$  consists of  $n$  data points from i.i.d. trajectories according to the behavior policy  $\pi_b$ . Then as  $n \rightarrow \infty$ ,*

$$\tilde{J}_n(\pi_\epsilon) \xrightarrow{\text{a.s.}} J(\pi_\epsilon).$$

The theorem implies that  $\tilde{J}(\pi_\epsilon)$  is not only a conservative estimator for  $J(\pi_\epsilon)$ , but under certain regularity conditions, it is also a consistent estimator of the policy's true value. Meaning that for large  $n$ , the trade-off between reliability and accuracy vanishes. In Appendix B.3 we prove this theorem and give the exact rate at which the HAMBO estimate converges to the true value of  $\pi_\epsilon$ . This rate depends on the choice of kernel, time horizon  $T$ , and dimensions of the environment  $(d_a, d_s)$ . As an example, if  $k$  is a Linear or RBF kernel, with high probability  $|\tilde{J}_n(\pi_\epsilon) - J(\pi_\epsilon)| = \tilde{\mathcal{O}}(n^{-1/2})$  where  $\tilde{\mathcal{O}}$  omits polylogarithmic factors.

The kernel assumptions in Theorem 3.8 hold for many popular kernels such as any inner product kernel  $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}')$  with Lipschitz continuous finite-dimensional feature maps  $\phi(\cdot)$  or smooth kernels such as the RBF. On the other hand, Theorem 3.8 does hold for non-smooth functions, e.g., those corresponding to a Matern since their maximum information capacity  $\gamma_n$  grows polynomially with  $n$ .

To the best of our knowledge, Theorem 3.8 is the first result that shows the consistency of a model-based finite-horizon OPE method for a continuous environments.

## 4 HAMBO WITH NEURAL NETWORKS

In practice, we often want to evaluate policies in settings where the state and action spaces are higher-dimensional, and have access to larger amounts of offline data. In such environments, GPs become unpractical as they tend to generalize poorly in high-dimensional domains and their inference becomes prohibitively expensive for larger datasets.

**The NN-based Statistical Model.** In this section, we discuss practical variants of HAMBO which employ neural networks that scale more favorably to large datasets and high-dimensional domains. Crucially, we need to be



able to quantify epistemic uncertainty. For this purpose, we employ Bayesian Neural Networks (BNNs) which model  $\mathbf{h}_\theta(\mathbf{s}, \mathbf{a})$  as a neural network function where  $\theta$  are the parameters of the neural network. BNNs presume a prior distribution  $p(\theta) = \mathcal{N}(\theta; 0, \lambda \mathbf{I})$  and maintain an approximation of the posterior  $p(\theta|\mathcal{D}_b) \propto p(\mathcal{D}_b|\theta)p(\theta)$  over neural network parameters. We use an independent Gaussian likelihood  $p(\mathcal{D}_b|\theta) = \prod_{i=1}^n \mathcal{N}(s'_i; \mathbf{h}_\theta(\mathbf{s}_i, \mathbf{a}_i), \boldsymbol{\nu}_\theta^2(\mathbf{s}_i, \mathbf{a}_i))$  where  $\boldsymbol{\nu}_\theta^2(\mathbf{s}, \mathbf{a})$  is the vector of transition noise variances which is also predicted by the BNN.

We use Stein Variational Gradient Descent (SVGD) [Liu and Wang, 2016] to approximate the posterior as a set of  $K$  particles  $\Theta = \{\theta_1, \dots, \theta_K\}$ . We form the mean prediction of our model as the average prediction of the  $K$  NNs:

$$\boldsymbol{\mu}_\Theta(\mathbf{s}, \mathbf{a}) = \frac{1}{K} \sum_{k=1}^K \mathbf{h}_{\theta_k}(\mathbf{s}, \mathbf{a}).$$

Similarly, we estimate the epistemic variance as

$$\boldsymbol{\sigma}_{\Theta, \epsilon}^2(\mathbf{s}, \mathbf{a}) = \frac{1}{K} \sum_{k=1}^K (\mathbf{h}_{\theta_k}(\mathbf{s}, \mathbf{a}) - \boldsymbol{\mu}_\Theta(\mathbf{s}, \mathbf{a}))^2.$$

The overall predictive distribution is the equally weighted mixture of all  $K$  NN-based conditional Gaussians, i.e.,

$$p(\mathbf{s}'|\mathbf{s}, \mathbf{a}, \mathcal{D}_b) = \frac{1}{K} \sum_{k=1}^K \mathcal{N}(\mathbf{s}'; \mathbf{h}_{\theta_k}(\mathbf{s}, \mathbf{a}), \boldsymbol{\nu}_{\theta_k}^2(\mathbf{s}, \mathbf{a})) \quad (5)$$

whose variance is  $\boldsymbol{\sigma}_\Theta^2(\mathbf{s}, \mathbf{a}) = \boldsymbol{\sigma}_{\Theta, \epsilon}^2(\mathbf{s}, \mathbf{a}) + \boldsymbol{\sigma}_{\Theta, a}^2(\mathbf{s}, \mathbf{a})$ , where  $\boldsymbol{\sigma}_{\Theta, a}^2(\mathbf{s}, \mathbf{a}) := \frac{1}{K} \sum_{k=1}^K \boldsymbol{\nu}_{\theta_k}^2(\mathbf{s}, \mathbf{a})$  represents aleatoric and  $\boldsymbol{\sigma}_{\Theta, \epsilon}^2(\mathbf{s}, \mathbf{a})$  the epistemic uncertainty.

**Calibrating the Model.** Since our BNN model uses approximate inference and a potentially misspecified prior, it may not satisfy the calibration condition of Assumption 3.1. Thus, we re-calibrate the model’s uncertainty estimates with a calibration set  $\mathcal{D}_c \subset \mathcal{D}_b$  that is withheld from the training. In particular, we use temperature scaling which chooses  $\tau > 0$  such that the scaled predictive distribution (5) with variance  $\tau^2 \boldsymbol{\sigma}_\Theta^2(\mathbf{s}, \mathbf{a})$  has a minimal empirical calibration error on  $\mathcal{D}_c$  [Kuleshov et al., 2018]. Algorithm 3 formalizes this technique. Note that re-calibrating the BNN model does not guarantee formal calibration in the sense of Assumption 3.1. However, in our experiments, we found it to reliably yield a conservative value estimate  $\tilde{J}(\pi_e)$ .

#### 4.1 PRACTICAL NN-BASED HAMBO VARIANTS

In the following, we discuss three ways of constructing adversarially hallucinated transition models based on our BNN model described in Equation (5). The formal pseudocode of all algorithms is presented in Appendix A.

**Continuous Adversary (HAMBO-CA).** This approach directly reflects the hallucinated adversarial transition model, introduced in (1) and (2). The adversary

$\eta(\mathbf{s}, \mathbf{a}) \in [-1, 1]^{d_s}$  chooses the mean of the Gaussian transition probability from the epistemic confidence set, i.e.,

$$\tilde{p}_\eta(\mathbf{s}'|\mathbf{s}, \mathbf{a}) := \mathcal{N}(\mathbf{s}'; \boldsymbol{\mu}_\Theta(\mathbf{s}, \mathbf{a}) + \tau^2 \eta(\mathbf{s}, \mathbf{a}) \boldsymbol{\sigma}_{\Theta, \epsilon}^2, \boldsymbol{\sigma}_{\Theta, a}^2(\mathbf{s}, \mathbf{a})).$$

To get the corresponding conservative value estimate  $\tilde{J}(\pi_e)$ , we need to solve the minimization problem  $\min_\eta J_{\tilde{p}_\eta}(\pi_e)$ . For this, we parameterize the adversary  $\eta(\mathbf{s}, \mathbf{a})$  as a neural network policy and use Soft Actor-Critic (SAC) [Haarnoja et al., 2018b] to maximize the negative return.

**Discrete Adversary (HAMBO-DA).** Our BNN posterior is approximated by a set of  $K$  NNs whose mean squared error difference corresponds to epistemic uncertainty. Thus, we can also construct a pessimistic transition model by letting the adversary choose which of the  $K$  NNs to pick. In this case, the adversary  $\vartheta(k|\mathbf{s}, \mathbf{a})$  is a categorical distribution over the NN indices  $\{1, \dots, K\}$ . The hallucinated transition model follows as:

$$\tilde{p}_\vartheta(\mathbf{s}'|\mathbf{s}, \mathbf{a}) := \sum_{k=1}^K \vartheta(k|\mathbf{s}, \mathbf{a}) \mathcal{N}(\mathbf{s}'; \mathbf{h}_{\theta_k}(\mathbf{s}, \mathbf{a}), \boldsymbol{\nu}_{\theta_k}^2(\mathbf{s}, \mathbf{a})).$$

Here, the adversary stochastically picks one of NN models at every step  $t = 0, \dots, T - 1$ . For this reason, we refer to this variant as DA1 (Algorithm 5). The corresponding value estimate follows as  $\tilde{J}_{\text{DA1}}(\pi_e) = \min_\vartheta J_{\tilde{p}_\vartheta}(\pi_e)$ . We solve the optimization problem by parameterizing the adversary  $\vartheta$  as a NN policy and use the clipped double DQN algorithm Fujimoto et al. [2018] to maximize the negative return.

Alternatively, we can constrain the adversary so that it has to commit to one of the  $K$  NN models for the entire trajectory. We refer to this variant as DAINF (Algorithm 6). In this case, the transition model corresponds to the predictive distribution of one of the NNs  $p_{\theta_k}(\mathbf{s}'|\mathbf{s}, \mathbf{a}) = \mathcal{N}(\mathbf{s}'; \mathbf{h}_{\theta_k}(\mathbf{s}, \mathbf{a}), \boldsymbol{\nu}_{\theta_k}^2(\mathbf{s}, \mathbf{a}))$ , and the value estimate follows as the minimum the policy values under each of the models, i.e.,  $\tilde{J}_{\text{DAInf}}(\pi_e) = \min_{k \in \{1, \dots, K\}} J_{p_{\theta_k}}(\pi_e)$ . If  $K$  is larger (e.g.,  $K > 20$ ), we recommend taking the empirical  $\delta$  quantile of the policy values  $\{J_{\tilde{p}_k}(\pi_e)\}_{k=1}^K$  instead of the minimum. In this case, DAINF has similarities to the model-based bootstrap approach of Kostrikov and Nachum [2020].

Naturally, the value estimates of DAINF are less pessimistic than those of DA1, i.e.  $\tilde{J}_{\text{DAInf}}(\pi_e) \geq \tilde{J}_{\text{DA1}}(\pi_e)$ , because the adversary cannot change which model it picks throughout the trajectory. In the experiment section, we investigate whether DAINF is still conservative enough to reliably yield lower bounds on the true policy values  $J(\pi_e)$ .

## 5 EXPERIMENTS

We start this section by illustrating the inner workings of HAMBO with a toy example to show why pessimism is crucial for COPE. We demonstrate that the convergence guarantees from Section 3.2 materialize in practice for

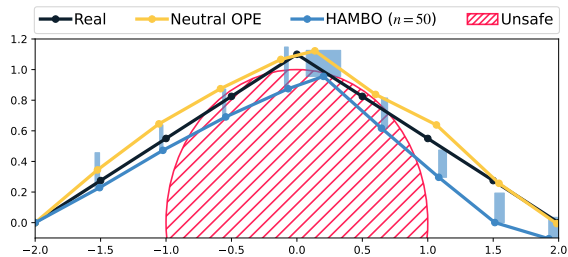


Figure 1: Hallucinated trajectories for model-based OPE and pessimistic HAMBO. While OPE overestimates the performance of the unsafe policy, HAMBO correctly gives a conservative estimates through its adversarial transition model. The adversary chooses the worst-case trajectory with the confidence sets (shaded blue areas).

GP models. Finally, we empirically evaluate and compare the practical variants of HAMBO with BNNs on various continuous control tasks. For comparability between our environments, we shift and scale all our results so that the true policy return value  $J(\pi_e)$  is 1.

## 5.1 ILLUSTRATIVE EXAMPLE

To illuminate the core idea of HAMBO and why pessimism is crucial for COPE, we conduct experiments on a toy environment which we call PointSafety (see Figure 1). In this environment, the agent navigates in the two-dimensional plane by applying actions  $\mathbf{a} \in [-0.5, 0.5]^2$  such that its position (i.e. state  $\mathbf{s} \in \mathcal{S} = \mathbb{R}^2$ ) changes to  $\mathbf{s}_{t+1} = \mathbf{s}_t + \mathbf{a}_t$ . The agent always starts on the left  $\mathbf{s}_0 = (-2, 0)$  and aims to go to its goal on the right  $\mathbf{s}_{\text{fin}} = (2, 0)$ . However, the unit circle is a danger zone, in which the agent is subject to highly negative rewards (red shaded area).

We consider evaluation policies  $\pi_y$  with an intermediate goal  $\mathbf{s}_{\text{im}} = (0, y)$  on the y-axis that goes in a straight line from  $\mathbf{s}_0$  to  $\mathbf{s}_{\text{im}}$  and then in a straight line from  $\mathbf{s}_{\text{im}}$  to the goal  $\mathbf{s}_{\text{fin}}$ . Note that policies  $\pi_y$  with  $|y| \leq 1.155$  are unsafe.

We generate an offline dataset by rolling out the behavior policy  $\pi_{1.6}$  with Gaussian action noise with a standard deviation of 0.1. Then, we evaluate  $\pi_{1.1}$ , which is unsafe (see black trajectory), by rolling it out using HAMBO-CA.

We compare this to a neutral variant that predicts the next state with the predictive mean  $\mu_{\Theta}(\mathbf{s}, \mathbf{a})$ , i.e., without pessimism. As we can observe from the yellow trajectory, it falsely estimates  $\pi_{1.1}$  as safe, that is, it predicts that the trajectory lies outside of the danger zone. The trajectories with the adversarial transition model and the corresponding epistemic confidence sets for every step are depicted in Fig 1. The adversary successfully moves the prediction towards the danger zone within the confidence set, and, thus, correctly estimates the policy to be unsafe. Overall, this demonstrates

a failure case of (neutral) off-policy evaluation and shows how HAMBO reliably gives a conservative estimate of the policy value through its pessimistic transition model.

## 5.2 EMPIRICAL CONVERGENCE OF HAMBO

For GP models, we show that HAMBO estimates converge to the true policy values (Theorem 3.8). Now, we empirically evaluate the behavior of GP-based HAMBO with an RBF kernel, as the number of offline data points grows. To this end, we consider two environments; a simple 2D PointEnv ( $\mathcal{S} = \mathbb{R}^2$ ,  $\mathcal{A} = [-1, 1]^2$ ), similar to the PointSafety environment, and the Pendulum-v1 environment from the OpenAI Gym [Brockman et al., 2016]. In the PointEnv, the agent has to navigate the origin and accordingly receives the negative distance to the origin as a reward.

To generate the offline dataset, we collect transition data by uniformly sampling states and actions from the state and action space respectively. For the PointEnv, we restrict the sampled states to  $[-40, 40]^2$  which covers the relevant part of the state space. As the evaluation policy, we use a proportional controller for the PointEnv, and a controller learned with SAC for the Pendulum.

Figure 2 plots the HAMBO estimates  $\tilde{J}(\pi_e)$  for a varying number of offline datapoints  $n = |\mathcal{D}_b|$ . We notice that in the PointEnv, when we have insufficient data (here, ca.  $n \leq 150$ ), the epistemic confidence regions of our GP model are large enough so that the transition model adversary sometimes manages to steer the policy outside the data support where the epistemic uncertainty is even higher. As a result, we see that  $\tilde{J}(\pi_e)$  are initially far below the true expected return  $J(\pi_e)$ . However, as  $n$  increases, the GP uncertainty regions become smaller, and, as we can observe in Figure 2,  $\tilde{J}(\pi_e)$  becomes an increasingly tighter lower bound, approaching  $J(\pi_e)$  for both the environments.

## 5.3 HAMBO FOR CONTINUOUS CONTROL

We evaluate the NN-based HAMBO methods from Section 4 on the continuous control tasks Pendulum-v1, Hopper-v3, and HalfCheetah-v3 from the OpenAI Gym and compare them to respective neutral (non-pessimistic) OPE methods.

Our general methodology is as follows: For a given environment, we first train a policy using the SAC algorithm Haarnoja et al. [2018a,b] and save several checkpoints of the agent. Then, some of the mediocre-performing checkpoints are rolled out to generate an offline dataset. After that, a given policy (usually one of the best checkpoints) is evaluated with the NN-based HAMBO variants.

We compare our approach to neutral OPE variants that do not use a pessimistic transition model [Fonteneau et al., 2013]. In particular, we consider various trajectory uncertainty propagation methods from Chua et al. [2018], em-

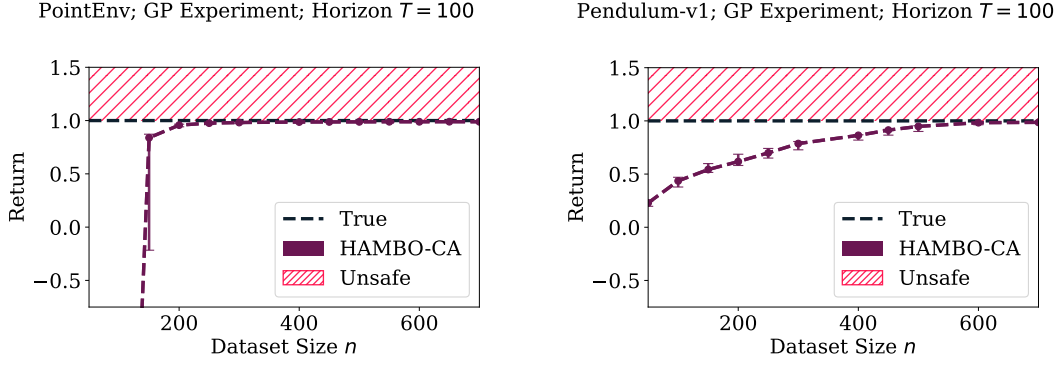


Figure 2: GP-based HAMBO for increasing offline dataset sizes  $n$  evaluated on the PointEnv and Pendulum-v1. The lower bound approaches the true return.

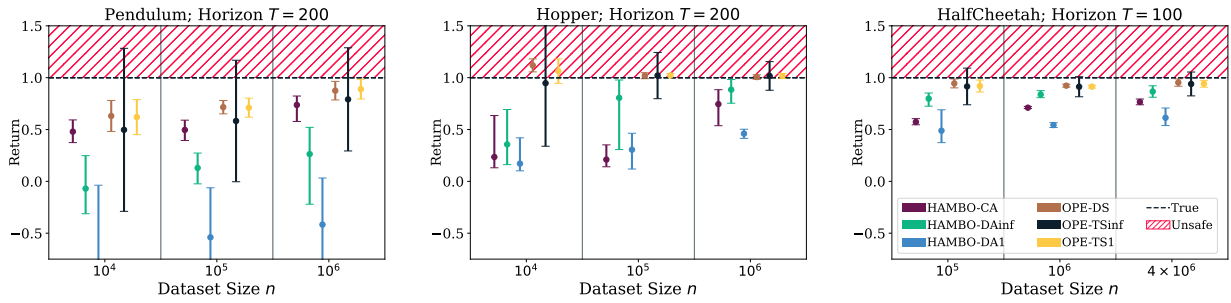


Figure 3: HAMBO variants and neutral OPE baselines for continuous control. Unlike neutral OPE, which frequently overestimates the true expected return, HAMBO always yields a valid lower bound, which becomes more accurate with  $n$ .

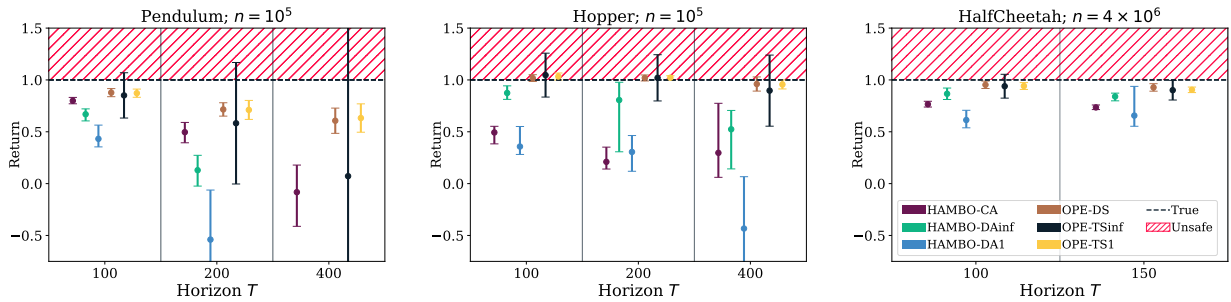


Figure 4: HAMBO variants and neutral OPE baselines for different horizons. With longer horizons, the variance of the neutral OPE estimates increases and HAMBO lower bounds become looser.

ployed in the context of OPE: First, we consider OPE-DS, where the transition model is approximated by a Gaussian  $p(s'|s, \mathbf{a}) = \mathcal{N}(s'; \mu_{\Theta}(s, \mathbf{a}); \sigma_{\Theta}^2(s, \mathbf{a}))$ , here the variance is the sum of the epistemic and aleatoric variance. Second, we consider OPE-TS1 where the transition model is the mixture of predictive Gaussians in (5). This means that, in every step, one of the NN models is chosen uniformly at random to compute the next state distribution. Third, we consider OPE-TSINF, where, for every episode, we randomly commit to one of the  $K$  NNs.

We investigate the following three aspects: 1) whether

a method yields reliable lower bounds, 2) the effect of the offline dataset size, and 3) the curse of long horizons. Figure 3 and 4 report the estimated expected policy returns, averaged over 5 seeds, alongside the corresponding confidence intervals.

**Reliable Lower Bounds.** The HAMBO variants are designed to give reliable lower bounds on the true expected return. The results in Figure 3 and 4 empirically confirm that, across all seeds, all NN-based HAMBO variants reliably provide lower bounds on  $J(\pi_e)$ , and, thus, fulfill the COPE requirements from Definition 2.1. In contrast, the neutral

OPE variants which do not introduce pessimism w.r.t. the epistemic uncertainty of the transition model fail to do so. In many cases, they overestimate the true policy value, particularly in the Hopper environment. This demonstrates the importance of pessimism in model-based COPE and affirms the validity of HAMBO, even with BNN models, where calibration (Definition 3.1) cannot be formally proven.

**Offline Dataset Size and Tightness.** The difference between HAMBO estimates  $\tilde{J}(\pi_e)$  and the true expected reward  $J(\pi_e)$  depends on the strength of the transition adversary, which is limited by the size of the epistemic confidence sets. As the size of the offline datasets  $\mathcal{D}_b$  increases, we can generally expect the epistemic uncertainty to shrink. Thus, the adversary  $\eta$  becomes less powerful and the HAMBO estimates become an increasingly tight lower bound.

In Figure 3, we empirically investigate this effect by varying the offline dataset size  $n$ . As we hypothesized, we can observe the general trend that the HAMBO estimates come close to the true policy value, as  $n$  increases. Moreover, we observe that the HAMBO-DA1 estimates are always strictly smaller than those of the HAMBO-DAINF variant. This is expected, since in the DA1 variant, the adversary can pick the worst-case NN transition model at every step while in the case of DAINF the adversary can only do so per trajectory, and, thus has less power. Since our experiment results indicate that the pessimism in HAMBO-DAINF is sufficient to obtain reliable lower bounds in practice, we conclude that HAMBO-DAINF is the preferred choice among the two. While HAMBO-DAINF performs better in Hopper and HalfCheetah, HAMBO-CA yields the tightest lower bounds in the Pendulum environment.

**The Curse of the Long Horizons.** Finally, we investigate the effect of the horizon length  $T$  on our COPE estimates. Over the course of a trajectory, the transition model estimation errors can compound and lead to large discrepancies. This is a well-studied phenomenon in model-based RL [e.g. see Janner et al., 2019]. In our case, this is reflected by the worst-case lower bound in Theorem 3.5 which depends exponentially on  $T$ .

To evaluate the empirical effect of horizon length, we report the (C)OPE estimates for an offline dataset of size  $n = 10^5$  across varying horizon lengths:  $T = 100, 200$  and  $400$  for the Pendulum and Hopper. For HalfCheetah, we only report horizon lengths of  $T = 100$  and  $150$ . Figure 4 displays the corresponding results. For an increasing horizon length, the variance of the neutral variants increases and the lower bounds of the conservative HAMBO estimates become looser. However, the observed decline in tightness in Figure 4 is much less pronounced than the exponential decline of the worst-case bound in Theorem 3.5.

For large horizon lengths, it can happen that the hallucinated trajectory under the pessimistic transition model strays far outside the support of the offline data. In such cases, unlike

neutral OPE methods, HAMBO will still provide lower bounds on the true expected return. However, these bounds can be very pessimistic. For instance, this can be observed in the case of Pendulum, where for  $T = 400$  the estimates of HAMBO-DA1 and HAMBO-DAINF go out of the chart. Making accurate long-horizon predictions is generally very hard. For instance, this is discussed extensively in the context of model-based RL in Janner et al. [2019]. Often, a discount factor is used when computing returns to alleviate these issues. We highlight that we work with undiscounted returns and continuous state-action spaces, and, thus, operate in the most challenging setting for OPE.

## 6 RELATED WORK

This work mainly contributes to the literature on off-policy evaluation for MDPs, which we divide to three categories.

**Model-Free OPE.** The key challenge in OPE is to the distribution shift between behavior and evaluated policy. A popular natural approach to correct the distribution mismatch is to use importance sampling (IS) ratios to re-weight the rewards collected by the behavior policy [Precup et al., 2000, Dudík et al., 2011] or to adjust the recursive updates when estimating the values directly via the Bellman equation [Precup et al., 2001, Sutton et al., 2015, Hallak and Mannor, 2017]. Some work also combine both approaches to obtain a more favorable bias-variance trade-off [Jiang and Li, 2016, Thomas and Brunskill, 2016]. Unlike HAMBO, these approaches are model-free, i.e., they do not learn a model of the state transitions. However, they suffer from three key disadvantages: First, they have notoriously high variance, especially if the evaluated policy differs a lot from the behavior policy Levine et al. [2020]. Second, they require the support of the behavior occupancy measure  $\rho^{\pi_b}$  to contain the support of  $\rho^{\pi_e}$  which is often not the case. In contrast, HAMBO still provides valid estimates in this scenario. Third, to compute the importance ratios, they assume access to the distribution of behavior policy which is almost never the case in practical applications where data is often collected by human experts. HAMBO does not require access to the behavior policy and, thus, is much more broadly applicable.

A recent line of work [Nachum et al., 2019a,b, Zhang et al., 2020, Yang et al., 2020] estimates the state occupancy correction ratios via a form of fixed point iteration, and does not require access to the behavior policy. However, the Bellmann-like fixed point iteration is not applicable to the finite horizon case that we study in this paper. In addition, due to the fixed point iteration, it is very hard to quantify the uncertainty or bound error that is associated with such OPE estimates, making them poorly suited to COPE.

**Model-Based OPE.** This approach first learns the transition dynamics, to then simulate rollouts with the evaluation policy  $\pi_e$  and thereby estimate the expected reward of  $\pi_e$  [e.g.,



Fonteneau et al., 2013, Hanna et al., 2017, Kostrikov and Nachum, 2020]. Due to error in predicting the transitions, the resulting OPE estimate may overestimate the policy’s performance which is prohibited in safety-critical applications. Our approach additionally simulates pessimistic trajectories using the model’s epistemic uncertainty, to avoid overestimation. Further, to the best of our knowledge, Theorem 3.8 is the first consistency result for model-based OPE.

**COPE and High-Confidence OPE.** We study the problem of COPE which seeks a high-probability lower bound on the expected return. This is closely related to estimating confidence bounds for OPE. Thomas et al. [2015] provide such confidence bounds for IS-based OPE estimates. However, due to the high variance of IS estimates, such bounds are often very loose [Levine et al., 2020]. Kallus and Uehara [2020] and Shi et al. [2021] propose a model-free approach to give asymptotically normal confidence intervals for directly  $J(\pi)$ . Assuming that the  $Q$ -function resides in an RKHS, Feng et al. [2020] and Feng et al. [2021] present rates of convergence, under theoretically unverified assumptions about the MDP. These model-free approaches only work for discounted, infinite-horizon MDPs, thus, are not generally applicable to our finite-horizon setting.

Hanna et al. [2017], Kostrikov and Nachum [2020] use model-based bootstrapping to construct confidence intervals for the OPE estimates. Kostrikov and Nachum [2020] prove the asymptotic correctness of the bootstrap confidence intervals only for finite state-action spaces. In contrast, we show the validity of our COPE estimates non-asymptotically for any  $|\mathcal{D}_b|$ , and in continuous state-action spaces. Alternatively, Fonteneau et al. [2009] and Paduraru [2013] employ a Lipschitz argument to obtain valid COPE estimates. Our derivation of the worst-case lower bound in Theorem 3.5 also uses Lipschitz continuity. However, the HAMBO estimate provide a tighter lower bound on the true policy value, as we use the local confidence intervals rather than the global Lipschitz constants to introduce pessimism. Furthermore, unlike the mentioned work, HAMBO does not require knowledge of the Lipschitz constant and works with sub-Gaussian noise.

## 7 CONCLUSION

HAMBO, a novel approach for COPE that forms a pessimistic estimate of the expected return by hallucinating adversarial trajectories within the epistemic confidence regions of the estimated transition model. We formally prove the validity and consistency of the resulting COPE estimates. We propose various scalable NN-based variants of HAMBO and empirically demonstrate that they give reliable and tight lower bounds on the true expected return.

Importantly, our approach does not require access to the probability distribution of the behavior policy and gives

reliable estimates, even when the support evaluation policy’s occupancy measure is not contained in the offline data distribution. This makes HAMBO particularly relevant for safety-critical real-world applications, where the offline data is mostly collected by human experts and we need to make reliable decisions about whether a given policy is good enough to be deployed.

HAMBO can be naturally combined with other offline reinforcement learning (ORL) algorithms to solve safety-critical ORL tasks. We leave this for future work to investigate.

## Acknowledgements

This research was supported by the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation program grant agreement no. 815943 and the Swiss National Science Foundation under NCCR Automation, grant agreement 51NF40 180545. Jonas Rothfuss was supported by an Apple Scholars in AI/ML fellowship. We thank Sebastian Curi for contributing to the initial idea for this project.

## References

- Felix Berkenkamp, Matteo Turchetta, Angela Schoellig, and Andreas Krause. Safe model-based reinforcement learning with stability guarantees. *Advances in Neural Information Processing Systems (NeurIPS)*, 30, 2017.
- Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. In *International Conference on Machine Learning (ICML)*, 2015.
- Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv*, 2016.
- Lukas Brunke, Melissa Greeff, Adam W Hall, Zhaocong Yuan, Siqi Zhou, Jacopo Panerati, and Angela P Schoellig. Safe learning in robotics: From learning-based control to safe reinforcement learning. *Annual Review of Control, Robotics, and Autonomous Systems*, 2022.
- Sayak Ray Chowdhury and Aditya Gopalan. On kernelized multi-armed bandits. In *International Conference on Machine Learning*, 2017.
- Kurtland Chua, Roberto Calandra, Rowan McAllister, and Sergey Levine. Deep reinforcement learning in a handful of trials using probabilistic dynamics models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- Sebastian Curi, Felix Berkenkamp, and Andreas Krause. Efficient model-based reinforcement learning through

- optimistic policy search and planning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- Miroslav Dudík, John Langford, and Lihong Li. Doubly robust policy evaluation and learning. In *International Conference on Machine Learning (ICML)*, 2011.
- Yihao Feng, Tongzheng Ren, Ziyang Tang, and Qiang Liu. Accountable off-policy evaluation with kernel bellman statistics. In *International Conference on Machine Learning (ICML)*, 2020.
- Yihao Feng, Ziyang Tang, Na Zhang, and Qiang Liu. Non-asymptotic confidence intervals of off-policy evaluation: Primal and dual bounds. In *International Conference on Learning Representations (ICLR)*, 2021.
- Raphael Fonteneau, Susan Murphy, Louis Wehenkel, and Damien Ernst. Inferring bounds on the performance of a control policy from a sample of trajectories. In *IEEE Symposium on Adaptive Dynamic Programming and Reinforcement Learning*, 2009.
- Raphael Fonteneau, Susan A Murphy, Louis Wehenkel, and Damien Ernst. Batch mode reinforcement learning based on the synthesis of artificial trajectories. *Annals of Operations Research*, 208:383–416, 2013.
- Scott Fujimoto, Herke van Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *International Conference on Machine Learning (ICML)*, 2018.
- Tuomas Haarnoja, Aurick Zhou, P. Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning (ICML)*, 2018a.
- Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, and Sergey Levine. Soft actor-critic algorithms and applications. *arXiv*, 2018b.
- Assaf Hallak and Shie Mannor. Consistent on-line off-policy evaluation. In *International Conference on Machine Learning (ICML)*, 2017.
- Josiah P. Hanna, Peter Stone, and Scott Niekum. Bootstrapping with models: Confidence intervals for off-policy evaluation. In *AAAI*, 2017.
- Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. When to trust your model: Model-based policy optimization. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- Nan Jiang and Lihong Li. Doubly robust off-policy value evaluation for reinforcement learning. In *International Conference on Machine Learning (ICML)*, 2016.
- Nathan Kallus and Masatoshi Uehara. Double reinforcement learning for efficient off-policy evaluation in markov decision processes. *J. Mach. Learn. Res.*, 2020.
- B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. *IEEE Transactions on Intelligent Transportation Systems*, 2021.
- Ilya Kostrikov and Ofir Nachum. Statistical bootstrapping for uncertainty estimation in off-policy evaluation. *arXiv*, 2020.
- Volodymyr Kuleshov, Nathan Fenner, and Stefano Ermon. Accurate uncertainties for deep learning using calibrated regression. In *International Conference on Machine Learning (ICML)*, 2018.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv*, 2020.
- Qiang Liu and Dilin Wang. Stein variational gradient descent: A general purpose bayesian inference algorithm. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2016.
- Travis Mandel, Yun-En Liu, Sergey Levine, Emma Brunskill, and Zoran Popovic. Offline policy evaluation across representations with applications to educational games. In *AAMAS*, 2014.
- Susan A Murphy, Mark J van der Laan, James M Robins, and Conduct Problems Prevention Research Group. Marginal mean models for dynamic regimes. *Journal of the American Statistical Association*, 2001.
- Ofir Nachum, Yinlam Chow, Bo Dai, and Lihong Li. Dualdice: Behavior-agnostic estimation of discounted stationary distribution corrections. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019a.
- Ofir Nachum, Bo Dai, Ilya Kostrikov, Yinlam Chow, Lihong Li, and Dale Schuurmans. Algaedice: Policy gradient from arbitrary experience. *arXiv*, 2019b.
- Cosmin Paduraru. *Off-policy evaluation in Markov decision processes*. PhD thesis, McGill University, 2013.
- Doina Precup, Richard S. Sutton, and Satinder P. Singh. Eligibility traces for off-policy policy evaluation. In *International Conference on Machine Learning (ICML)*, 2000.

- Doina Precup, Richard S Sutton, and Sanjoy Dasgupta. Off-policy temporal-difference learning with function approximation. In *International Conference on Machine Learning (ICML)*, 2001.
- Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning)*. The MIT Press, 2005.
- Chengchun Shi, Runzhe Wan, Victor Chernozhukov, and Rui Song. Deeply-debiased off-policy interval estimation. In *International Conference on Machine Learning (ICML)*, 2021.
- Richard S. Sutton, A. Rupam Mahmood, and Martha White. An emphatic approach to the problem of off-policy temporal-difference learning. *arXiv*, 2015.
- Philip Thomas and Emma Brunskill. Data-efficient off-policy policy evaluation for reinforcement learning. In *International Conference on Machine Learning (ICML)*, 2016.
- Philip Thomas, Georgios Theodorou, and Mohammad Ghavamzadeh. High-confidence off-policy evaluation. *AAAI*, 2015.
- Mengjiao Yang, Ofir Nachum, Bo Dai, Lihong Li, and Dale Schuurmans. Off-policy evaluation via the regularized lagrangian. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- Ruiyi Zhang, Bo Dai, Lihong Li, and Dale Schuurmans. Gendice: Generalized offline estimation of stationary values. In *International Conference on Learning Representations (ICLR)*, 2020.