

---

# Hallucinated Adversarial Control for Conservative Offline Policy Evaluation - Supplementary Material

---

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

## A ALGORITHM AND EXPERIMENT DETAILS

In the following, we provide algorithmic formalizations and implementation details of the HAMBO framework and its practical variants which were discussed in the main paper.

### A.1 GENERIC HAMBO ALGORITHM FROM SECTION 3.1

First, we formalize the general HAMBO framework from Section 3.1:

---

**Algorithm 1** HAMBO Framework

---

**Require:** Offline dataset  $\mathcal{D}_b$ , evaluation policy  $\pi_e$ , reward function  $r(\cdot, \cdot)$ , Horizon  $T$ , initial state distribution  $p_0(\mathbf{s}_0)$   
 $(\boldsymbol{\mu}_n, \boldsymbol{\sigma}_n, \beta_n) \leftarrow \text{TrainModel}(\mathcal{D}_b)$  ▷ Train statistical model with offline data  
 $\tilde{p}_\eta(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t) \leftarrow p_e(\mathbf{s}_{t+1} - \boldsymbol{\mu}_n(\mathbf{s}_t, \mathbf{a}_t) - \beta_n \boldsymbol{\eta}(\mathbf{s}_t, \mathbf{a}_t) \boldsymbol{\sigma}_n(\mathbf{s}_t, \mathbf{a}_t))$  ▷ Set up adversarial transition model.  
 $\tilde{J}(\pi) \leftarrow \min_{\boldsymbol{\eta}} \mathbb{E}_{\mathbf{s}_0 \sim p_0} [\mathbb{E}_{p_\eta, \pi} [\sum_{t=0}^T r(\mathbf{s}_t, \mathbf{a}_t)]]$  ▷ Optimize adversary to get pessimistic value estimate.  
**return**  $\tilde{J}(\pi)$

---

We estimate  $J_{\tilde{p}_\eta}(\pi_e) = \mathbb{E}_{\mathbf{s}_0 \sim p_0} [\mathbb{E}_{p_\eta, \pi} [\sum_{t=0}^T r(\mathbf{s}_t, \mathbf{a}_t)]]$  via Monte Carlo estimation, i.e., we roll out  $L$  trajectories and estimate the expectation as the average of the trajectory return:

$$\hat{J}_{\tilde{p}_\eta}(\pi_e) = \frac{1}{L} \sum_{l=1}^L \sum_{t=0}^T r(\mathbf{s}_{l,t}, \mathbf{a}_{l,t}) \text{ where } \mathbf{s}_{l,0} \sim p_0, \mathbf{a}_{l,t} \sim \pi(\mathbf{a}|\mathbf{s}_{l,t}), \mathbf{s}_{l,t+1} \sim \tilde{p}_\eta(\mathbf{s}'|\mathbf{s}_{l,t}, \mathbf{a}_{l,t}) \quad (1)$$

The optimization of the adversary corresponds to a standard optimal control problem for which we use traditional methods such as trajectory optimization or model-free RL algorithms such as SAC.

### A.2 BNN BASED HAMBO VARIANTS

#### A.2.1 The BNN model

We use fully connected neural networks with 4 hidden layers each of size 256 with ReLU activation functions. Before training, the offline data inputs and targets are standardized. The NN takes the concatenated state and action as input (i.e.,  $d_s + d_a$  dimensional) and outputs a vector of size  $2d_s$  which is split into two vectors of size  $d_s$ . The first one corresponds to

---

\*Equal contribution.

the mean prediction  $\mathbf{h}_\theta(\mathbf{s}, \mathbf{a})$  and the second one is the raw the aleatoric standard deviation which is fed through a softplus function to ensure positivity of  $\nu_{\theta}^2(\mathbf{s}, \mathbf{a})$

As BNN prior we use a standard Normal distribution over the NN parameters  $\theta$ , i.e.,  $p(\theta) = \mathcal{N}(\theta; \mathbf{0}, \mathbf{I})$ . However, as commonly done for BNNs to alleviate the problems of prior misspecification, we add a temperature parameter  $\tau$  to the prior, so that we have  $p(\theta|\mathcal{D}_b) \propto p(\mathcal{D}_b|\theta)p(\theta)^\tau$ . This hyper-parameter is chosen to as  $\tau = 0.0001$  for Pendulum and Hopper and  $\tau = 0.01$  for the HalfCheetah control environment.

We use Stein Variational Gradient Descent (SVGD) [Liu and Wang, 2016] for approximate posterior inference. In particular, we approximate the posterior  $p(\theta|\mathcal{D}_b)$  with  $K$  NN particles  $\{\theta_1, \dots, \theta_K\}$ . After randomly initializing the parameters of the  $K$  NNs, the parameters are iteratively updated with the SVGD update rule:

$$\theta_k \leftarrow \theta_k + \frac{1}{K} \sum_{k'=1}^K [k(\theta_{k'}, \theta_k) \nabla_{\theta_{k'}} \log p(\theta_{k'}|\mathcal{D}_b) + \nabla_{\theta_{k'}} k(\theta_{k'}, \theta_k)] \quad \forall k = 1, \dots, K. \quad (2)$$

Here  $k(\cdot, \cdot)$  is a kernel function on the space of NN parameters vectors. In our experiments, we use an RBF kernel  $k(\theta, \theta') = \exp\left(-\|\theta - \theta'\|_2^2 / (2\ell)\right)$  with a length scale of  $\ell = 10$  and  $K = 5$  NN particles. Note that the kernel here is different from the one in Section 3.2. Algorithm 2 summarizes how to obtain the SVGD BNN posterior approximation.

---

### Algorithm 2 SVGD

---

**Require:** Training data  $\mathcal{D}$ , number of particles  $K$

Initialize NN parameter vectors  $\theta_1, \dots, \theta_K$

▷ Initialize SVGD particles.

**while** not converged **do**

$\log p(\theta_k|\mathcal{D}) \leftarrow \log p(\mathcal{D}|\theta_k) + \tau \log p(\theta_k) \quad \forall k = 1, \dots, K$

$\theta_k \leftarrow \theta_k + \frac{1}{K} \sum_{k'=1}^K [k(\theta_{k'}, \theta_k) \nabla_{\theta_{k'}} \log p(\theta_{k'}|\mathcal{D}) + \nabla_{\theta_{k'}} k(\theta_{k'}, \theta_k)] \quad \forall k = 1, \dots, K$

**end while**

**return**  $\{\mathbf{h}_{\theta_1}, \nu_{\theta_1}^2, \dots, \mathbf{h}_{\theta_K}, \nu_{\theta_K}^2\}$

▷ Return the  $K$  NN predictive mean and aleatoric variance functions

---

### A.3 RECALIBRATION OF THE BNN UNCERTAINTY ESTIMATES

To obtain well-calibrated confidence sets for HAMBO, we recalibrate the BNNs predictive distribution. In particular, we use temperature scaling based on the regression calibration error Kuleshov et al. [2018]. We perform re-calibration based on the predictive distribution  $\mathcal{N}(\mu_\Theta(\mathbf{s}, \mathbf{a}), \sigma_\Theta^2(\mathbf{s}, \mathbf{a}))$ . The calibration error compares the predictive quantiles of this Normal distribution with the corresponding empirical frequencies of data points, that fall below the predicted quantiles. Formally, we define  $\Phi_\tau^{-1}(\alpha; \mathbf{s}, \mathbf{a}) : [0, 1]^{d_s} \rightarrow \mathcal{S}$  as the quantile function (inverse cumulative density function) of  $\mathcal{N}(\mu_\Theta(\mathbf{s}, \mathbf{a}), \tau^2 \sigma_\Theta^2(\mathbf{s}, \mathbf{a}))$  where  $\tau \in \mathbb{R}^{d_s}$  is the temperature scaling vector. Given a calibration dataset  $\mathcal{D}_c = \{(\mathbf{s}, \mathbf{a}, \mathbf{s}')\}$ , the calibration error [Kuleshov et al., 2018] for multivariate distributions follows as

$$\text{CalErr}(\tau) := \frac{1}{d_s} \sum_{j=1}^{d_s} \frac{1}{|A|} \sum_{\alpha \in A} (\text{EmpFreq}(\alpha, \tau)_j - \alpha)^2, \quad (3)$$

where  $A = \{0.1, \dots, 0.9, 0.99\}$  is a set of confidence levels and

$$\text{EmpFreq}(\alpha; \tau) := \frac{1}{|\mathcal{D}_c|} \sum_{(\mathbf{s}, \mathbf{a}, \mathbf{s}') \in \mathcal{D}_c} \mathbf{1}\{\mathbf{s}' \leq \Phi_\tau^{-1}(\alpha; \mathbf{s}, \mathbf{a})\} \quad (4)$$

is a vector-valued function of the (per dimension) empirical frequencies of the prediction targets that fall below the  $\alpha$  quantile. Finally, we recalibrate the BNN predictions, by choosing the variance scaling vector  $\tau$  such that the calibration error is minimized, i.e., we choose

$$\tau^* = \arg \min_{\tau} \text{CalErr}(\tau). \quad (5)$$

Algorithm 3 summarizes this BNN re-calibration procedure:

---

**Algorithm 3** CALIBRATEBNN

---

**Require:** calibration dataset  $\mathcal{D}_c$ , predictive mean  $\mu(\cdot, \cdot)$ , predictive variance  $\sigma^2(\cdot, \cdot)$   
     $A \leftarrow \{0.1, \dots, 0.9, 0.99\}$ . ▷ Fix a set of confidence levels.  
     $\Phi^{-1}(\cdot; \mathbf{s}, \mathbf{a}, \boldsymbol{\tau}) : [0, 1] \mapsto \mathbb{R}^{d_s}$  as the inverse CDF of the Gaussian distribution  $\mathcal{N}(\mu(\mathbf{s}, \mathbf{a}), \tau^2 \sigma^2(\mathbf{x}_i))$ .  
    Define  $\text{EmpFreq}(\alpha; \boldsymbol{\tau}) \leftarrow \frac{1}{|\mathcal{D}_c|} \sum_{(\mathbf{s}, \mathbf{a}, \mathbf{s}') \in \mathcal{D}_c} \mathbf{1}\{\mathbf{s}' \leq \Phi_{\boldsymbol{\tau}}^{-1}(\alpha; \mathbf{s}, \mathbf{a}, \boldsymbol{\tau})\}$   
    Define  $\text{CalErr}(\boldsymbol{\tau}) \leftarrow \frac{1}{d_s} \sum_{j=1}^{d_s} \frac{1}{|A|} \sum_{\alpha \in A} (\text{EmpFreq}(\alpha, \boldsymbol{\tau})_j - \alpha)^2$   
    **return**  $\arg \min_{\boldsymbol{\tau}} \text{CalErr}(\boldsymbol{\tau})$  ▷ Choose  $\boldsymbol{\tau}$  that minimizes the calibration error

---

#### A.4 THE NN-BASED HAMBO VARIANTS

Here, we provide algorithmic descriptions of the NN-Based HAMBO variants from Section 4 as well as details about their implementation and how the corresponding experiments were conducted.

##### A.4.1 HAMBO with a Continuous Adversary (HAMBO-CA)

HAMBO-CA directly reflects the hallucinated adversarial transition model, introduced in Section 3. The adversary  $\boldsymbol{\eta}(\mathbf{s}, \mathbf{a}) \in [-1, 1]^{d_s}$  chooses the mean of the Gaussian transition probability from the epistemic confidence set, i.e.,

$$\tilde{p}_{\boldsymbol{\eta}}(\mathbf{s}' | \mathbf{s}, \mathbf{a}) := \mathcal{N}(\mathbf{s}'; \mu_{\Theta}(\mathbf{s}, \mathbf{a}) + \tau^2 \boldsymbol{\eta}(\mathbf{s}, \mathbf{a}) \sigma_{\Theta, e}^2, \sigma_{\Theta, a}^2(\mathbf{s}, \mathbf{a})) \quad (6)$$

For obtaining the corresponding conservative value estimate  $\tilde{J}(\pi_e) = \min_{\boldsymbol{\eta}} J_{\tilde{p}_{\boldsymbol{\eta}}}(\pi_e)$ , we need to find the adversary  $\boldsymbol{\eta}^*$  that minimizes the expected return. For this, we parameterize the adversary  $\boldsymbol{\eta}$  as a neural network policy with two hidden layers of size 256 with ReLU activations and a tanh squashed Gaussian conditional distribution over the adversary actions in  $[-1, 1]^{d_s}$ . We use SAC [Haarnoja et al., 2018] to maximize the negative expected return of the adversary policy. As usual, to stabilize the SAC training and avoid Q-value overestimation, we use double critics and trailing target critics. The SAC training is conducted in rounds consisting of rollouts of 1000 episodes under the hallucinated transition model where actions are chosen by  $\pi_e$ , followed by 1000 gradient steps on the SAC objectives. For the gradient steps, we use a batch size of 1024 and the Adam optimizer with a learning rate of  $10^{-3}$  for critic and policy and  $5 * 10^{-5}$  for the SAC entropy parameter. After SAC has converged, we take the adversary policy  $\boldsymbol{\eta}^*$  and estimate the expected return  $\hat{J}_{\tilde{p}_{\boldsymbol{\eta}^*}}(\pi_e)$  of  $\pi_e$  under the adversary transition model, induced by  $\boldsymbol{\eta}^*$  with  $L = 10^4$  trajectories (see Eq. 1). The HAMBO-CA method is summarized in Algorithm 4.

---

**Algorithm 4** HAMBO-CA

---

**Require:** Offline dataset  $\mathcal{D}_b$ , evaluation policy  $\pi_e$ , Number of BNN particles  $K$   
    Select a subset of  $\mathcal{D}_b$  as calibration set  $\mathcal{D}_c$   
     $\{\mathbf{h}_{\theta_1}, \boldsymbol{\nu}_{\theta_1}^2, \dots, \mathbf{h}_{\theta_K}, \boldsymbol{\nu}_{\theta_K}^2\} \leftarrow \text{SVGD}(\mathcal{D}_b \setminus \mathcal{D}_c, K)$  ▷ Train BNN via SVGD and get predictive NN functions  
     $\mu_{\Theta}(\mathbf{s}, \mathbf{a}) \leftarrow \frac{1}{K} \sum_{k=1}^K \mathbf{h}_{\theta_k}(\mathbf{s}, \mathbf{a})$  ▷ Calculate posterior mean.  
     $\sigma_{\Theta, e}^2(\mathbf{s}, \mathbf{a}) \leftarrow \frac{1}{K} \sum_{k=1}^K (\mathbf{h}_{\theta_k}(\mathbf{s}, \mathbf{a}) - \mu_{\Theta}(\mathbf{s}, \mathbf{a}))^2$  ▷ Calculate epistemic uncertainty.  
     $\sigma_{\Theta, a}^2(\mathbf{s}, \mathbf{a}) \leftarrow \frac{1}{K} \sum_{k=1}^K \boldsymbol{\nu}_{\theta_k}^2(\mathbf{s}, \mathbf{a})$  ▷ Calculate aleatoric uncertainty.  
     $\tau \leftarrow \text{CalibrateBNN}(\mathcal{D}_c, \mu_{\Theta}, \sigma_{\Theta, e}^2 + \sigma_{\Theta, a}^2)$  ▷ Calibrate the model  
    Initialize adversary policy  $\boldsymbol{\eta}$   
     $\tilde{p}_{\boldsymbol{\eta}}(\mathbf{s}' | \mathbf{s}, \mathbf{a}) \leftarrow \mathcal{N}(\mathbf{s}'; \mu_{\Theta}(\mathbf{s}, \mathbf{a}) + \tau^2 \boldsymbol{\eta}(\mathbf{s}, \mathbf{a}) \sigma_{\Theta, e}^2, \sigma_{\Theta, a}^2(\mathbf{s}, \mathbf{a}))$  ▷ Setup hallucinated adversarial transition model  
     $\boldsymbol{\eta}^* \leftarrow \text{SoftActorCritic}(-J_{\tilde{p}_{\boldsymbol{\eta}}}(\pi_e), \boldsymbol{\eta})$  ▷ Train adversary  $\boldsymbol{\eta}$  via SAC to maximize the negative return  
     $\tilde{J}(\pi_e) \leftarrow \hat{J}_{\tilde{p}_{\boldsymbol{\eta}^*}}(\pi_e)$  ▷ Estimate expected return of  $\pi_e$  via sampling (see Eq. 1)  
    **return**  $\tilde{J}(\pi_e)$

---

##### A.4.2 HAMBO with a Discrete Adversary (HAMBO-DA1 and HAMBO-DAINF)

In the case of HAMBO-DA1 the adversary  $\vartheta$  has discrete action  $\{1, \dots, K\}$ , i.e. picking one of the  $K$  particles. We parameterize the adversary policy as a neural network with two hidden layers of size 256 with ReLU activations and softmax-categorical distribution over the  $K$  discrete actions. To train this adversary policy, we use clipped double DQN

[Fujimoto et al., 2018]. The double DQN training is conducted in rounds consisting of rollouts of 1000 episodes under the hallucinated transition model where actions are chosen by  $\pi_e$ , followed by 1000 gradient steps on the DQN objectives. For the gradient steps, we use a batch size of 1024 and the Adam optimizer with a learning rate of  $10^{-3}$ . Once double DQN has converged, we take the adversary policy  $\vartheta^*$  and estimate the expected return  $\hat{J}_{\tilde{p}_{\vartheta^*}}(\pi_e)$  of  $\pi_e$  under the adversary transition model, induced by  $\vartheta^*$  with  $L = 10^4$  trajectories (see Eq. 1). The overall HAMBO-DA1 method is summarized in Algorithm 5.

---

**Algorithm 5** HAMBO-DA1

---

**Require:** Offline dataset  $\mathcal{D}_b$ , evaluation policy  $\pi_e$ , Number of BNN particles  $K$   
 $\{\mathbf{h}_{\theta_1}, \boldsymbol{\nu}_{\theta_1}^2, \dots, \mathbf{h}_{\theta_K}, \boldsymbol{\nu}_{\theta_K}^2\} \leftarrow \text{SVGD}(\mathcal{D}_b, K)$  ▷ Train BNN via SVGD and get predictive NN functions  
Initialize adversary policy  $\vartheta$   
 $\tilde{p}_{\vartheta}(s'|\mathbf{s}, \mathbf{a}) := \sum_{k=1}^K \vartheta(k|\mathbf{s}, \mathbf{a}) \mathcal{N}(s'; \mathbf{h}_{\theta_k}(\mathbf{s}, \mathbf{a}), \boldsymbol{\nu}_{\theta_k}^2(\mathbf{s}, \mathbf{a}))$  ▷ Setup hallucinated adversarial transition model  
 $\vartheta^* \leftarrow \text{DoubleDQN}(-J_{\tilde{p}_{\vartheta}}(\pi_e), \vartheta)$  ▷ Train adversary  $\vartheta$  via to maximize the negative return  
 $\tilde{J}(\pi_e) \leftarrow \hat{J}_{\tilde{p}_{\vartheta^*}}(\pi_e)$  ▷ Estimate expected return of  $\pi_e$  via sampling (see Eq. 1)  
**return**  $\tilde{J}(\pi_e)$

---

In contrast to HAMBO-DA1, HAMBO-DAINF uses a weaker adversary that has to commit to one of the BNN particles for the entire trajectory. As a result, the corresponding pessimistic HAMBO estimate can simply be chosen as the minimum of the expected evaluation policy return under each of the NN models in the particle approximation, i.e.  $J(\pi_e) = \min_{k \in \{1, \dots, K\}} J_{p_{\theta_k}}(\pi_e)$ . The HAMBO-DAINF method is summarized in Algorithm 5.

---

**Algorithm 6** HAMBO-DAINF

---

**Require:** Offline dataset  $\mathcal{D}_b$ , evaluation policy  $\pi_e$ , Number of BNN particles  $K$   
 $\{\mathbf{h}_{\theta_1}, \boldsymbol{\nu}_{\theta_1}^2, \dots, \mathbf{h}_{\theta_K}, \boldsymbol{\nu}_{\theta_K}^2\} \leftarrow \text{SVGD}(\mathcal{D}_b, K)$  ▷ Train BNN via SVGD and get predictive NN functions  
 $p_{\theta_k}(s'|\mathbf{s}, \mathbf{a}) \leftarrow \mathcal{N}(s'; \mathbf{h}_{\theta_k}(\mathbf{s}, \mathbf{a}), \boldsymbol{\nu}_{\theta_k}^2(\mathbf{s}, \mathbf{a}))$   
 $\tilde{J}(\pi_e) \leftarrow \min_{k \in \{1, \dots, K\}} \hat{J}_{p_{\theta_k}}(\pi_e)$  ▷ Estimate return of  $\pi_e$  for each model (see Eq. 1) and take minimum  
**return**  $\tilde{J}(\pi_e)$

---

## B PROOFS AND DERIVATIONS

*Proof of Proposition 3.2.* By Assumption 3.1 we have, with probability  $1 - \delta$ , uniformly over  $\mathcal{S} \times \mathcal{A}$ , that

$$|\boldsymbol{\mu}_n(\mathbf{s}, \mathbf{a}) - f(\mathbf{s}, \mathbf{a})| \leq \beta_n(\delta) \boldsymbol{\sigma}_n(\mathbf{s}, \mathbf{a}). \quad (7)$$

Hence, there exists an (adversary) mapping  $\boldsymbol{\eta}^\dagger : \mathcal{S} \times \mathcal{A} \mapsto [-1, 1]^{d_s}$  such that every  $\forall \mathbf{s}, \mathbf{a} \in \mathcal{S} \times \mathcal{A}$  we have

$$f(\mathbf{s}, \mathbf{a}) = \boldsymbol{\mu}_n(\mathbf{s}, \mathbf{a}) + \beta_n(\delta) \boldsymbol{\eta}^\dagger(\mathbf{s}, \mathbf{a}) \boldsymbol{\sigma}_n(\mathbf{s}, \mathbf{a}), \quad (8)$$

and, thus the hallucinated transition model is equal to the true transition dynamics, i.e.,

$$\tilde{p}_{\boldsymbol{\eta}^\dagger}(s_{t+1}|\mathbf{s}_t, \mathbf{a}_t) = p_e(s_{t+1} - \boldsymbol{\mu}_n(\mathbf{s}, \mathbf{a}) + \beta \boldsymbol{\eta}(\mathbf{s}, \mathbf{a})^\dagger \boldsymbol{\sigma}_n(\mathbf{s}, \mathbf{a})) = p_e(s_{t+1} - f(\mathbf{s}, \mathbf{a})) = p(s_{t+1}|\mathbf{s}_t, \mathbf{a}_t). \quad (9)$$

Finally, we can use this to show

$$\tilde{J}(\pi_e) := \min_{\boldsymbol{\eta}} J_{\tilde{p}_{\boldsymbol{\eta}}}(\pi_e) \leq J_{\tilde{p}_{\boldsymbol{\eta}^\dagger}}(\pi_e) = J_p(\pi_e) = J(\pi_e), \quad (10)$$

which concludes the proof. □

*Proof of Example 3.4.* If  $\pi(\mathbf{a}|\mathbf{s})$  can be reparametrized as  $g(\mathbf{s}, \boldsymbol{\zeta})$ , where  $\boldsymbol{\zeta} \sim p(\boldsymbol{\zeta})$  and  $g$  is  $L_g$ -Lipschitz, we have that the two random variables are equal in distribution, i.e.  $\mathbf{a} \stackrel{d}{=} g(\mathbf{s}, \boldsymbol{\zeta})$ . Therefore,

$$\mathcal{W}_1(\pi(\mathbf{a}|\mathbf{s}), \pi(\mathbf{a}'|\mathbf{s}')) = \inf_{\gamma \in \Gamma(\pi(\mathbf{a}|\mathbf{s}), \pi(\mathbf{a}'|\mathbf{s}'))} \mathbb{E}_{\mathbf{a}, \mathbf{a}' \sim \gamma} [\|\mathbf{a} - \mathbf{a}'\|_2] \quad (11)$$

$$\leq \mathbb{E}_{\mathbf{a}, \mathbf{a}' \sim \tilde{\gamma}} [\|\mathbf{a} - \mathbf{a}'\|_2] = \mathbb{E}_{\boldsymbol{\zeta}} [\|g(\mathbf{s}, \boldsymbol{\zeta}) - g(\mathbf{s}', \boldsymbol{\zeta})\|_2] \quad (12)$$

$$\leq L_g \|\mathbf{s} - \mathbf{s}'\|_2 \quad (13)$$

where  $\tilde{\gamma}(\mathbf{a}, \mathbf{a}')$  is the joint probability distribution of  $(g(\mathbf{s}, \boldsymbol{\zeta}), g(\mathbf{s}', \boldsymbol{\zeta}))$ , and, thus a coupling. Hence, we have shown that  $\pi(\mathbf{a}, \mathbf{s})$  is  $L_g$ -Lipschitz w.r.t. the Wasserstein-1 distance. □

## B.1 PROOF OF THEOREM 3.5

The following lemmata will be used to prove the theorem.

**Lemma B.1** (Reparameterizability of two random variables with covariates). *Let  $X$  and  $Y$  be random variables with finite expectation and corresponding probability distributions  $p$  and  $q$ . Then, we can reparameterize  $Y$  as  $Y \stackrel{d}{=} X + \omega|_X$ , where  $\omega|_X$  is a covariate that is generally dependent on  $X$  and satisfies*

$$\mathbb{E}_X \mathbb{E}_{\omega|X} [\|\omega\|_2] = \mathcal{W}_1(p, q), \quad (14)$$

where  $\mathcal{W}_1(p, q)$  is the Wasserstein-1 distance between  $p$  and  $q$ .

*Proof.* Recall that the Wasserstein-1 distance is defined as infimum over couplings between  $p$  and  $q$ , i.e.,

$$\mathcal{W}_1 = \inf_{\gamma \in \Gamma(p, q)} \mathbb{E}_{X', Y' \sim \gamma} [\|X' - Y'\|_2] \quad (15)$$

If the expectation of  $p$  and  $q$  is finite, then the infimum over couplings in (15) is attained for some  $\gamma^*(x, y)$ . Now we construct the covariate  $\omega|_X$  which is defined by applying the change of variable  $g_x(x, y) \mapsto (x, y - x) = (x, \omega)$  to  $\gamma^*$ , so that we get  $\tilde{\gamma}^*(x, \omega) = \gamma^*(x, x + \omega)$ . The conditional distribution of the covariate  $\omega|_X$  is

$$\tilde{\gamma}^*(\omega|x) = \frac{\gamma^*(x, x + \omega)}{\gamma^*(x)} = \frac{\gamma^*(x, x + \omega)}{p(x)}$$

Now, given our construction of  $\omega|_X$ , we aim to show that  $Y \stackrel{d}{=} X + \omega|_X$ . Define the random variable  $Z := X + \omega|_X$ . Then we have

$$p(z) = \int_{\mathcal{X}} p(x, z - x) dx = \int_{\mathcal{X}} p(x) \underbrace{\tilde{\gamma}^*(z - x|x)}_{\omega} dx \quad (16)$$

$$= \int_{\mathcal{X}} \gamma^*(x, x + (z - x)) dx = \int_{\mathcal{X}} \gamma^*(x, z) dx = q(z) \quad (17)$$

which shows that the pdf of  $z$  is  $q$ , the probability density of  $Y$ . Since  $\gamma^*(x, y)$  is the coupling that minimizes the transport cost, we can write

$$\mathcal{W}_1(p, q) = \inf_{\gamma \in \Gamma(p, q)} \mathbb{E}_{x', y' \sim \gamma} [\|x' - y'\|_2] = \mathbb{E}_{x', y' \sim \gamma^*} [\|x' - y'\|_2] = \mathbb{E}_{x' \sim p(x')} \mathbb{E}_{\omega \sim \tilde{\gamma}^*(\omega|x')} [\|\omega\|_2] \quad (18)$$

which shows that  $\mathbb{E}_X \mathbb{E}_{\omega|X} [\|\omega\|_2] = \mathcal{W}_1(p, q)$ , and, thus concludes the proof.  $\square$

**Corollary B.2.** *Let  $\pi(\mathbf{a}|s)$  be  $L_\pi$ -Lipschitz w.r.t. the Wasserstein-1 distance. For any arbitrary but fixed  $s, s' \in \mathcal{S}$  we denote  $A$  and  $A'$  as the random variables that follow the conditional distributions  $\pi(\mathbf{a}|s)$  and  $\pi(\mathbf{a}'|s')$  respectively. Then, we can construct a covariate  $\omega|_A$  such that  $A' \stackrel{d}{=} A + \omega|_A$  and  $\mathbb{E}_A \mathbb{E}_{\omega|A} [\|\omega\|_2] \leq L_\pi \|s - s'\|_2$ .*

*Proof.* The corollary directly follows from Lemma B.1 and the definition of the  $L_\pi$ -Lipschitz continuity w.r.t. the Wasserstein-1, i.e., that  $\forall s, s' \in \mathcal{S}$  we have that  $\mathcal{W}_1(\pi(\mathbf{a}|s), \pi(\mathbf{a}'|s')) \leq L_\pi \|s - s'\|_2$ .  $\square$

**Lemma B.3** (Lipschitz continuity of Wasserstein-one distance implies Lipschitz continuity in expectation). *Let  $f : \mathcal{X}_1 \times \mathcal{X}_2 \rightarrow \mathcal{Y}$  be  $L_f$  Lipschitz continuous and  $x_2$  a random variable with distribution  $p(\cdot|x_1)$  that is  $L_p$  Lipschitz w.r.t. the Wasserstein-1 distance. Then we have*

$$\mathbb{E}_{x_2 \sim p(\cdot|x_1)} [f(x_1, x_2)] - \mathbb{E}_{x_2 \sim p(\cdot|x_1')} [f(x_1', x_2')] \leq \bar{L}_f \|x_1 - x_1'\|.$$

with  $\bar{L}_f = L_f(1 + L_p)$ .

*Proof.*

$$\begin{aligned}
\mathbb{E}_{x_2 \sim p(\cdot|x_1)} [f(x_1, x_2)] - \mathbb{E}_{x'_2 \sim p(\cdot|x'_1)} [f(x'_1, x'_2)] &= \mathbb{E}_{x_2 \sim p(\cdot|x_1)} [f(x_1, x_2)] - \mathbb{E}_{x_2 \sim p(\cdot|x_1)} \left[ \mathbb{E}_{\omega \sim \tilde{\gamma}^*(\omega|x_2)} [f(x'_1, x_2 + \omega)] \right] \\
&\quad \text{(Lemma B.1)} \\
&= \mathbb{E}_{x_2 \sim p(\cdot|x_1)} \left[ \mathbb{E}_{\omega \sim \tilde{\gamma}^*(\omega|x_2)} [f(x_1, x_2) - f(x'_1, x_2 + \omega)] \right] \\
&\leq L_f \mathbb{E}_{x_2 \sim p(\cdot|x_1)} \left[ \mathbb{E}_{\omega \sim \tilde{\gamma}^*(\omega|x_2)} [\|x_1 - x'_1\|_2 + \|\omega\|_2] \right] \quad \text{(Lipschitzness of } f) \\
&\leq L_f \|x_1 - x'_1\|_2 + L_f L_p \|x_1 - x'_1\|_2 \quad \text{(Corollary B.2)} \\
&= L_f(1 + L_p) \|x_1 - x'_1\|_2.
\end{aligned}$$

□

In the following, we bound the difference between the pessimistic and true return with the distance between the true and pessimistic trajectory using the Lipschitz continuity of the reward function and the policy's Wasserstein-one distance.

**Lemma B.4** (Bound on difference between pessimistic and true return estimate). *Under Assumption 3.3 we have*

$$\left| J(\pi_e) - \tilde{J}(\pi_e) \right| \leq \bar{L}_r \mathbb{E}_{\epsilon_{0:T-1}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} \left[ \sum_{t=0}^{T-1} \|\mathbf{s}_t - \tilde{\mathbf{s}}_t\|_2 \right] \right].$$

where  $\bar{L}_r = L_r(1 + L_\pi)$ .

*Proof.* We have

$$\begin{aligned}
\left| J(\pi_e) - \tilde{J}(\pi_e) \right| &= \left| \mathbb{E}_{\mathbf{s}_{0:T}, \mathbf{a}_{0:T}} \left[ \sum_{t=0}^{T-1} r(\mathbf{s}_t, \mathbf{a}_t) \right] - \mathbb{E}_{\tilde{\mathbf{s}}_{0:T}, \tilde{\mathbf{a}}_{0:T}} \left[ \sum_{t=0}^{T-1} r(\tilde{\mathbf{s}}_t, \tilde{\mathbf{a}}_t) \right] \right| \\
&= \left| \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \sum_{t=0}^{T-1} r(\mathbf{s}_t, \mathbf{a}_t) \right] - \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} \left[ \sum_{t=0}^{T-1} r(\tilde{\mathbf{s}}_t, \mathbf{a}_t + \omega_t) \right] \right] \right| \quad \text{(Lemma B.1)} \\
&= \left| \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} \left[ \sum_{t=0}^{T-1} r(\mathbf{s}_t, \mathbf{a}_t) - r(\tilde{\mathbf{s}}_t, \mathbf{a}_t + \omega_t) \right] \right] \right|.
\end{aligned}$$

From Lemma B.1, know that  $\mathbb{E} \omega = \mathcal{W}_1(\pi(\cdot|\mathbf{s}_t), \pi(\cdot|\tilde{\mathbf{s}}_t))$ , where  $\pi(\cdot|\cdot)$  is continuous w.r.t. the WD-1 distance. Therefore,

$$\begin{aligned}
\left| J(\pi_e) - \tilde{J}(\pi_e) \right| &\leq \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} \left[ \sum_{t=0}^{T-1} L_r(1 + L_\pi) \|\mathbf{s}_t - \tilde{\mathbf{s}}_t\|_2 \right] \right] \quad \text{(Lemma B.3)} \\
&= \bar{L}_r \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} \left[ \sum_{t=0}^{T-1} \|\mathbf{s}_t - \tilde{\mathbf{s}}_t\|_2 \right] \right]. \quad (\bar{L}_r = L_r(1 + L_\pi))
\end{aligned}$$

□

Next, we bound the distance between the true and pessimistic trajectory with the epistemic uncertainty around the true trajectory.

**Lemma B.5** (Bound on pessimistic and true trajectory). *Under Assumption 3.1 and 3.3 with probability at least  $1 - \delta$  for all  $\eta : \mathcal{S} \rightarrow [-1, 1]^{d_s}$  we have for all  $t \in \{0, \dots, T\}$  that*

$$\mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|\mathbf{s}_{t+1} - \tilde{\mathbf{s}}_{t+1}\|_2] \right] \leq (1 + \sqrt{d_s}) \beta \sum_{i=0}^{(t+1)-1} (\bar{L}_f + (1 + \sqrt{d_s}) \beta \bar{L}_\sigma)^{(t+1)-1-i} \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|\boldsymbol{\sigma}_n(\mathbf{s}_i, \mathbf{a}_i)\|_2] \right].$$

*Proof.* We prove by induction. For  $t = 1$  we have

$$\begin{aligned}
\mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|\mathbf{s}_1 - \tilde{\mathbf{s}}_1\|_2] \right] &= \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|f(\mathbf{s}_0, \mathbf{a}_0) + \epsilon_0 - \boldsymbol{\mu}_n(\mathbf{s}_0, \mathbf{a}_0) - \beta \boldsymbol{\sigma}_n(\mathbf{s}_0, \mathbf{a}_0) \boldsymbol{\eta}(\mathbf{s}_0, \mathbf{a}_0) - \epsilon_0\|_2] \right] \\
&\hspace{15em} \text{(Lemma B.1)} \\
&\leq \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|f(\mathbf{s}_0, \mathbf{a}_0) - \boldsymbol{\mu}_n(\mathbf{s}_0, \mathbf{a}_0)\|_2 + \|\beta \boldsymbol{\sigma}_n(\mathbf{s}_0, \mathbf{a}_0) \boldsymbol{\eta}(\mathbf{s}_0)\|_2] \right] \\
&\leq \left(1 + \sqrt{d_s}\right) \beta_n \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|\boldsymbol{\sigma}_n(\mathbf{s}_0, \mathbf{a}_0)\|_2] \right] \hspace{5em} (\boldsymbol{\eta} \in [-1, 1]^{d_s})
\end{aligned}$$

We get the induction hypothesis that for an arbitrary but fixed  $t \geq 0$  we have

$$\mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|\mathbf{s}_t - \tilde{\mathbf{s}}_t\|_2] \right] \leq (1 + \sqrt{d_s}) \beta_n \sum_{i=0}^{t-1} (\bar{L}_f + (1 + \sqrt{d_s}) \beta \bar{L}_\sigma)^{t-1-i} \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|\boldsymbol{\sigma}_n(\mathbf{s}_i, \mathbf{a}_i)\|_2] \right]$$

Now for the induction step we can first derive

$$\begin{aligned}
\mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|\mathbf{s}_{t+1} - \tilde{\mathbf{s}}_{t+1}\|_2] \right] &= \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|f(\mathbf{s}_t, \mathbf{a}_t) + \epsilon_t - \boldsymbol{\mu}_n(\tilde{\mathbf{s}}_t, \tilde{\mathbf{a}}_t) - \beta_n \boldsymbol{\sigma}_n(\tilde{\mathbf{s}}_t, \tilde{\mathbf{a}}_t) \boldsymbol{\eta}(\tilde{\mathbf{s}}_t, \tilde{\mathbf{a}}_t) - \epsilon_t\|_2] \right] \\
&= \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|f(\mathbf{s}_t, \mathbf{a}_t) - \boldsymbol{\mu}_n(\tilde{\mathbf{s}}_t, \mathbf{a}_t + \boldsymbol{\omega}_t) - \beta_n \boldsymbol{\sigma}_n(\tilde{\mathbf{s}}_t, \mathbf{a}_t + \boldsymbol{\omega}_t) \boldsymbol{\eta}(\tilde{\mathbf{s}}_t, \mathbf{a}_t + \boldsymbol{\omega}_t)\|_2] \right] \\
&\hspace{15em} \text{(Lemma B.1)} \\
&\leq \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|f(\mathbf{s}_t, \mathbf{a}_t) - f(\tilde{\mathbf{s}}_t, \mathbf{a}_t + \boldsymbol{\omega}_t) + f(\tilde{\mathbf{s}}_t, \mathbf{a}_t + \boldsymbol{\omega}_t) - \boldsymbol{\mu}_n(\tilde{\mathbf{s}}_t, \mathbf{a}_t + \boldsymbol{\omega}_t)\|_2] \right] \\
&\quad + \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|\beta_n \boldsymbol{\sigma}_n(\tilde{\mathbf{s}}_t, \mathbf{a}_t + \boldsymbol{\omega}_t) \boldsymbol{\eta}(\tilde{\mathbf{s}}_t, \mathbf{a}_t + \boldsymbol{\omega}_t)\|_2] \right] \\
&\leq \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|f(\mathbf{s}_t, \mathbf{a}_t) - f(\tilde{\mathbf{s}}_t, \mathbf{a}_t + \boldsymbol{\omega}_t)\|_2 + \|f(\tilde{\mathbf{s}}_t, \mathbf{a}_t + \boldsymbol{\omega}_t) - \boldsymbol{\mu}_n(\tilde{\mathbf{s}}_t, \mathbf{a}_t + \boldsymbol{\omega}_t)\|_2] \right] \\
&\quad + \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|\beta_n \boldsymbol{\sigma}_n(\tilde{\mathbf{s}}_t, \mathbf{a}_t + \boldsymbol{\omega}_t) \boldsymbol{\eta}(\tilde{\mathbf{s}}_t, \mathbf{a}_t + \boldsymbol{\omega}_t)\|_2] \right] \\
&\leq \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\bar{L}_f \|\mathbf{s}_t - \tilde{\mathbf{s}}_t\|_2 + (1 + \sqrt{d_s}) \beta_n \|\boldsymbol{\sigma}_n(\tilde{\mathbf{s}}_t, \mathbf{a}_t + \boldsymbol{\omega}_t)\|_2] \right] \\
&\hspace{15em} \text{(Lemma B.3 and } \boldsymbol{\eta} \in [-1, 1]^{d_s})
\end{aligned}$$

By applying the triangle inequality and adding and subtracting  $\boldsymbol{\sigma}_n$  to the second term,

$$\begin{aligned}
\mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|\mathbf{s}_{t+1} - \tilde{\mathbf{s}}_{t+1}\|_2] \right] &\leq \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\bar{L}_f \|\mathbf{s}_t - \tilde{\mathbf{s}}_t\|_2 + (1 + \sqrt{d_s}) \beta_n \|\boldsymbol{\sigma}_n(\tilde{\mathbf{s}}_t, \mathbf{a}_t + \boldsymbol{\omega}_t)\|_2] \right] \\
&= \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\bar{L}_f \|\mathbf{s}_t - \tilde{\mathbf{s}}_t\|_2] \right] \\
&\quad + \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [(1 + \sqrt{d_s}) \beta_n \|\boldsymbol{\sigma}_n(\tilde{\mathbf{s}}_t, \mathbf{a}_t + \boldsymbol{\omega}_t) - \boldsymbol{\sigma}_n(\mathbf{s}_t, \mathbf{a}_t) + \boldsymbol{\sigma}_n(\mathbf{s}_t, \mathbf{a}_t)\|_2] \right] \\
&\leq \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\bar{L}_f \|\mathbf{s}_t - \tilde{\mathbf{s}}_t\|_2 + (1 + \sqrt{d_s}) \beta_n (\|\boldsymbol{\sigma}_n(\tilde{\mathbf{s}}_t, \mathbf{a}_t + \boldsymbol{\omega}_t) - \boldsymbol{\sigma}_n(\mathbf{s}_t, \mathbf{a}_t)\|_2)] \right] \\
&\quad + (1 + \sqrt{d_s}) \beta_n \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|\boldsymbol{\sigma}_n(\mathbf{s}_t, \mathbf{a}_t)\|_2] \right] \\
&\leq \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\bar{L}_f \|\mathbf{s}_t - \tilde{\mathbf{s}}_t\|_2 + (1 + \sqrt{d_s}) \beta_n (\bar{L}_\sigma \|\mathbf{s}_t - \tilde{\mathbf{s}}_t\|_2 + \|\boldsymbol{\sigma}_n(\mathbf{s}_t, \mathbf{a}_t)\|_2)] \right] \\
&\hspace{15em} \text{(Lemma B.1)} \\
&= \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [(\bar{L}_f + (1 + \sqrt{d_s}) \beta_n \bar{L}_\sigma) \|\mathbf{s}_t - \tilde{\mathbf{s}}_t\|_2 + (1 + \sqrt{d_s}) \beta_n \|\boldsymbol{\sigma}_n(\mathbf{s}_t, \mathbf{a}_t)\|_2] \right]
\end{aligned}$$

Next, we apply the induction hypothesis

$$\begin{aligned}
\mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|\mathbf{s}_{t+1} - \tilde{\mathbf{s}}_{t+1}\|_2] \right] &\leq \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} \left[ \left( \bar{L}_f + (1 + \sqrt{d_s}) \beta_n \bar{L}_\sigma \right) \|\mathbf{s}_t - \tilde{\mathbf{s}}_t\|_2 + (1 + \sqrt{d_s}) \beta_n \|\boldsymbol{\sigma}_n(\mathbf{s}_t, \mathbf{a}_t)\|_2 \right] \right] \\
&\leq \left[ \left( \bar{L}_f + (1 + \sqrt{d_s}) \beta_n \bar{L}_\sigma \right) (1 + \sqrt{d_s}) \beta_n \right] \\
&\quad \times \left( \sum_{i=0}^{t-1} \left( \bar{L}_f + (1 + \sqrt{d_s}) \beta_n \bar{L}_\sigma \right)^{t-1-i} \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|\boldsymbol{\sigma}_n(\mathbf{s}_i, \mathbf{a}_i)\|_2] \right] \right. \\
&\quad \left. + \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|\boldsymbol{\sigma}_n(\mathbf{s}_t, \mathbf{a}_t)\|_2] \right] \right) \\
&= (1 + \sqrt{d_s}) \beta_n \sum_{i=0}^{(t+1)-1} \left( \bar{L}_f + (1 + \sqrt{d_s}) \beta_n \bar{L}_\sigma \right)^{(t+1)-1-i} \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|\boldsymbol{\sigma}_n(\mathbf{s}_i, \mathbf{a}_i)\|_2] \right]
\end{aligned}$$

□

Using the above lemmas, we present the proof to the main theorem.

**Proof of Theorem 3.5.**

$$\begin{aligned}
|J(\pi_e) - \tilde{J}(\pi_e)| &\leq \bar{L}_r \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} \left[ \sum_{t=0}^{T-1} \|\mathbf{s}_t - \tilde{\mathbf{s}}_t\|_2 \right] \right] && \text{(Lemma B.4)} \\
&\leq \bar{L}_r \sum_{t=0}^{T-1} (1 + \sqrt{d_s}) \beta_n \sum_{i=0}^{(t+1)-1} \left( \bar{L}_f + (1 + \sqrt{d_s}) \beta_n \bar{L}_\sigma \right)^{(t+1)-1-i} \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} [\|\boldsymbol{\sigma}_n(\mathbf{s}_i, \mathbf{a}_i)\|_2] \right] \\
&&& \text{(Lemma B.5)}
\end{aligned}$$

Since  $\bar{L}_f \geq 1$ , then for all  $0 \leq i \leq t$  and  $0 \leq t \leq T-1$ ,

$$\left( \bar{L}_f + (1 + \sqrt{d_s}) \beta_n \bar{L}_\sigma \right)^{(t+1)-1-i} \leq \left( \bar{L}_f + (1 + \sqrt{d_s}) \beta_n \bar{L}_\sigma \right)^{T-1}$$

which allows us to write,

$$\begin{aligned}
|J(\pi_e) - \tilde{J}(\pi_e)| &\leq \bar{L}_r (1 + \sqrt{d_s}) \beta_n \left( 1 + \bar{L}_f + (1 + \sqrt{d_s}) \beta_n \bar{L}_\sigma \right)^{T-1} T \mathbb{E}_{\epsilon_{0:T}, \mathbf{a}_{0:T}} \left[ \mathbb{E}_{\omega_{0:T}} \left[ \sum_{t=0}^{T-1} \|\boldsymbol{\sigma}_n(\mathbf{s}_t, \mathbf{a}_t)\|_2 \right] \right] \\
&= \bar{L}_r (1 + \sqrt{d_s}) \beta_n \left( 1 + \bar{L}_f + (1 + \sqrt{d_s}) \beta_n \bar{L}_\sigma \right)^{T-1} T \mathbb{E}_{\mathbf{s}_{0:T}, \mathbf{a}_{0:T}} \left[ \sum_{t=0}^{T-1} \|\boldsymbol{\sigma}_n(\mathbf{s}_t, \mathbf{a}_t)\|_2 \right] \\
&= \left[ \bar{L}_r (1 + \sqrt{d_s}) \beta_n \left( 1 + \bar{L}_f + (1 + \sqrt{d_s}) \beta_n \bar{L}_\sigma \right)^{T-1} T \right] \\
&= \left[ \bar{L}_r (1 + \sqrt{d_s}) \beta_n \left( 1 + \bar{L}_f + (1 + \sqrt{d_s}) \beta_n \bar{L}_\sigma \right)^{T-1} T \right] \\
&\quad \times \left( \int_{\mathcal{S} \times \mathcal{A}} \sum_{t=0}^{T-1} p(\mathbf{s}_t = \mathbf{s}, \mathbf{a}_t = \mathbf{a} | \pi_e, \mathcal{M}) \|\boldsymbol{\sigma}_n(\mathbf{s}, \mathbf{a})\|_2 d\mathbf{s} d\mathbf{a} \right) \\
&= \bar{L}_r (1 + \sqrt{d_s}) \beta_n \left( 1 + \bar{L}_f + (1 + \sqrt{d_s}) \beta_n \bar{L}_\sigma \right)^{T-1} T \int_{\mathcal{S} \times \mathcal{A}} T \rho^{\pi_e}(\mathbf{s}, \mathbf{a}) \|\boldsymbol{\sigma}_n(\mathbf{s}, \mathbf{a})\|_2 d\mathbf{s} d\mathbf{a} \quad \text{(See Eq. 2)} \\
&= \bar{L}_r (1 + \sqrt{d_s}) \beta_n \left( 1 + \bar{L}_f + (1 + \sqrt{d_s}) \beta_n \bar{L}_\sigma \right)^{T-1} T^2 \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \rho^{\pi_e}} [\|\boldsymbol{\sigma}_n(\mathbf{s}, \mathbf{a})\|_2]
\end{aligned}$$

□

In summary, the deviation between the true and pessimistic return is proportional to the expected epistemic uncertainty of the evaluation policy state-occupancy measure  $\rho^{\pi_e}$ , and the constant  $C_n$  defined as

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

In Appendix B.3, we provide consistency guarantees for our method. In particular, we prove under further assumptions on the true dynamics function  $f$ , that  $\left|J(\pi_e) - \tilde{J}(\pi_e)\right| \rightarrow 0$ , for  $n \rightarrow \infty$ .

## B.2 PROOF OF KNOWN RESULTS FOR KERNEL METHODS

We first recall the notion of maximum mutual information [Srinivas et al., 2012, Cover and Thomas, 2006]. The mutual information  $I(\mathbf{x}_{1:n}; k)$  quantifies the reduction in uncertainty due to the observations  $\mathbf{x}_{1:n}$ . Given a GP model  $\mathcal{GP}(0, k(\cdot, \cdot))$  and gaussian noise assumption, mutual information is equal to

$$I(\mathbf{x}_{1:n}) = \frac{1}{2} \log \det(\mathbf{I} + \sigma_\epsilon^{-2} \mathbf{K})$$

with the kernel matrix  $\mathbf{K} = [k(\mathbf{x}_i, \mathbf{x}_j)]_{i,j \leq n}$ . The maximum information capacity or maximum mutual information of a kernel  $k$  is an upper bound on the mutual information, and is defined as

$$\gamma_n = \max_{\mathbf{x}_{1:n}} I(\mathbf{x}_{1:n}).$$

Table 1 shows the growth rate of  $\gamma_n$  with  $n$  for multiple different kernels.

*Proof of Lemma 3.6.* Let  $\gamma_n$  be the maximum mutual information of  $\mathcal{GP}(0, k(\cdot, \cdot))$ . Set  $\beta_n(\delta) := (B + \sigma_\epsilon \sqrt{2(\gamma_n + 1 + \ln(d_s/\delta))})$ . Element-wise application of Theorem 2 in Chowdhury and Gopalan [2017] over the dimensions of  $\mathcal{S}$  and taking a union bound proves the lemma.  $\square$

*Proof of Lemma 3.7.* First, we prove the Lipschitz continuity of  $\mathbf{f}$ . By the Cauchy-Schwartz inequality, we have  $\forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}$

$$|f_j(\mathbf{x}) - f_j(\mathbf{x}')| = |\langle f_j, k(\mathbf{x}, \cdot) - k(\mathbf{x}', \cdot) \rangle_k| \leq \|f_j\|_k d_k(\mathbf{x}, \mathbf{x}') \quad (19)$$

Since  $\|f_j\|_k \leq B, \forall j = 1, \dots, d_s$  and  $d_k(\mathbf{x}, \mathbf{x}')$  is  $L_k$ -Lipschitz, we have that

$$\|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{x}')\|_2 = \sqrt{\sum_{j=1}^{d_s} (f_j(\mathbf{x}) - f_j(\mathbf{x}'))^2} \leq \sqrt{d_s B^2 d_k(\mathbf{x}, \mathbf{x}')^2} = \sqrt{d_s} B d_k(\mathbf{x}, \mathbf{x}') \leq \sqrt{d_s} B L_k \|\mathbf{x} - \mathbf{x}'\|_2; \quad (20)$$

Next, we show the Lipschitz continuity of the GP standard deviation. By Lemma 12 in Curi et al. [2020], we have, independent of  $n$ ,  $|\sigma_n(\mathbf{x}) - \sigma_n(\mathbf{x}')| \leq d_k(\mathbf{x}, \mathbf{x}')$  for the GP standard deviation. Now, we make a similar argument as above:

$$\|\sigma_n(\mathbf{x}) - \sigma_n(\mathbf{x}')\|_2 \leq \sqrt{d_s d_k^2(\mathbf{x}, \mathbf{x}')^2} \leq \sqrt{d_s} L_k \|\mathbf{x} - \mathbf{x}'\|_2 \quad (21)$$

which shows that  $\sigma_n(\cdot)$  is  $\sqrt{d_s} L_k$ -Lipschitz.  $\square$

## B.3 PROOF OF THEOREM 3.8

For showing consistency of our lower bound in Theorem 3.5 for the GP case, we first prove that the uncertainty with respect to an i.i.d., data sampling distribution  $p(x)$  shrinks in expectation.

**Lemma B.6** (Shrinking uncertainty in expectation). *Let  $p(x)$  denote a data sampling distribution with compact support. Then the following holds for sequences  $\{x_i\}_{i=0}^{n-1}$  sampled i.i.d. from  $p(x)$ ,*

$$C_n^2 \mathbb{E}_{\mathbf{x}_{1:n} \sim p} [\sigma^2(x_n | \{x_i\}_{i=0}^{n-1})] \leq C_n^2 \mathbb{E}_{\mathbf{x}_{1:n-1} \sim p} [\sigma^2(x_{n-1} | \{x_i\}_{i=0}^{n-2})]. \quad (22)$$

*Proof.*

$$\begin{aligned}
C_n^2 \mathbb{E}_{\mathbf{x}_{1:n} \sim p} [\sigma^2(x_n | \{x_i\}_{i=0}^{n-1})] &= C_n^2 \mathbb{E}_{x_{1:n-1} \sim p} \left[ \mathbb{E}_{x_n \sim p} [\sigma^2(x_n | \{x_i\}_{i=0}^{n-1})] \right] \\
&= C_n^2 \mathbb{E}_{x_{1:n-1} \sim p} \left[ \mathbb{E}_{x \sim p} [\sigma^2(x | \{x_i\}_{i=0}^{n-1})] \right] \\
&\leq C_n^2 \mathbb{E}_{x_{1:n-1} \sim p} \left[ \mathbb{E}_{x \sim p} [\sigma^2(x | \{x_i\}_{i=0}^{n-2})] \right] && \text{(Monotonicity of variance)} \\
&= C_n^2 \mathbb{E}_{x_{1:n-2} \sim p} \left[ \mathbb{E}_{x_{n-1} \sim p} \left[ \mathbb{E}_{x \sim p} [\sigma^2(x | \{x_i\}_{i=0}^{n-2}) | x_{n-1}] \right] \right] \\
&= C_n^2 \mathbb{E}_{x_{1:n-2} \sim p} \left[ \mathbb{E}_{x \sim p} [\sigma^2(x | \{x_i\}_{i=0}^{n-2})] \right] && \text{(All points are sampled i.i.d from } p) \\
&= C_n^2 \mathbb{E}_{x_{1:n-2} \sim p} \left[ \mathbb{E}_{x_{n-1} \sim p} [\sigma^2(x_{n-1} | \{x_i\}_{i=0}^{n-2})] \right] \\
&= C_n^2 \mathbb{E}_{x_{1:n-1} \sim p} [\sigma^2(x_{n-1} | \{x_i\}_{i=0}^{n-2})].
\end{aligned}$$

□

**Lemma B.7** (Bound on expectation of uncertainty at  $n$ ). *Let  $p(\mathbf{x})$  denote the data sampling distribution with a compact support. Then the following holds for sequences  $\{\mathbf{x}_i\}_{i=0}^{n-1}$  sampled i.i.d. from  $p(\mathbf{x})$ ,*

$$nC_n^2 \mathbb{E}_{\mathbf{x}_{1:n} \sim p} [\sigma^2(\mathbf{x}_n | \{\mathbf{x}_i\}_{i=0}^{n-1})] \leq C_n^2 \mathbb{E}_{\mathbf{x}_{1:n} \sim p} \left[ \sum_{j=1}^n \sigma^2(\mathbf{x}_j | \{\mathbf{x}_i\}_{i=0}^{j-1}) \right]. \quad (23)$$

Moreover, we have

$$C_n^2 \mathbb{E}_{\mathbf{x}_{1:n} \sim p} [\sigma^2(\mathbf{x}_n | \{\mathbf{x}_i\}_{i=0}^{n-1})] \leq \frac{C_n^2 \gamma_n}{n},$$

where  $\gamma_n$  represents the maximum information gain (Srinivas et al. [2012], Cover and Thomas [2006]).

*Proof.* We prove by induction. For  $n = 1$ , Eq. 23 holds trivially. Now assume  $n > 1$ ,

$$\begin{aligned}
nC_n^2 \mathbb{E}_{\mathbf{x}_{1:n} \sim p} [\sigma^2(\mathbf{x}_n | \{\mathbf{x}_i\}_{i=0}^{n-1})] &= C_n^2 \mathbb{E}_{\mathbf{x}_{1:n} \sim p} [\sigma^2(\mathbf{x}_n | \{\mathbf{x}_i\}_{i=0}^{n-1})] + (n-1)C_n^2 \mathbb{E}_{\mathbf{x}_{1:n} \sim p} [\sigma^2(\mathbf{x}_n | \{\mathbf{x}_i\}_{i=0}^{n-1})] \\
&\leq C_n^2 \mathbb{E}_{\mathbf{x}_{1:n} \sim p} [\sigma^2(\mathbf{x}_n | \{\mathbf{x}_i\}_{i=0}^{n-1})] + (n-1)C_n^2 \mathbb{E}_{\mathbf{x}_{1:n-1} \sim p} [\sigma^2(\mathbf{x}_{n-1} | \{\mathbf{x}_i\}_{i=0}^{n-2})] \\
&&& \text{(Lemma B.6)} \\
&\leq C_n^2 \mathbb{E}_{\mathbf{x}_{1:n} \sim p} [\sigma^2(\mathbf{x}_n | \{\mathbf{x}_i\}_{i=0}^{n-1})] + C_n^2 \mathbb{E}_{\mathbf{x}_{1:n-1} \sim p} \left[ \sum_{j=1}^{n-1} \sigma^2(\mathbf{x}_j | \{\mathbf{x}_i\}_{i=0}^{j-1}) \right] \\
&&& \text{(By induction hypothesis)} \\
&= C_n^2 \mathbb{E}_{\mathbf{x}_{1:n} \sim p} \left[ \sum_{j=1}^n \sigma^2(\mathbf{x}_j | \{\mathbf{x}_i\}_{i=0}^{j-1}) \right].
\end{aligned}$$

Note,  $\sum_{j=1}^n \sigma^2(\mathbf{x}_{j+1} | \{\mathbf{x}_i\}_{i=0}^j)$  is a measure of the mutual information associated to the sampling scheme, and lower bounds the mutual information. The mutual information  $I(\mathbf{x}_{1:n})$  quantifies the reduction in uncertainty due to the observations  $\mathbf{x}_{1:n}$  Cover and Thomas [2006]. When  $f \in \mathcal{H}_k$ , mutual information is equal to

$$I(\mathbf{x}_{1:n}) = \frac{1}{2} \log \det(\mathbf{I} + \lambda^{-1} \mathbf{K})$$

with the kernel matrix  $\mathbf{K} = [k(\mathbf{x}_i, \mathbf{x}_j)]_{i,j \leq n}$ . Moreover,

$$\mathbb{E}_{\mathbf{x}_{1:n} \sim p} \left[ \sum_{j=1}^n \sigma^2(x_{j+1} | \{x_i\}_{i=0}^j) \right] \leq I(\mathbf{x}_{1:n}).$$

The maximum information gain, is an upper bound on the mutual information, and is defined as

$$\gamma_n = \max_{\mathbf{x}_{1:n}} I(\mathbf{x}_{1:n}).$$

Therefore, by definition of  $\gamma_n$ , it is greater than the mutual information of all sampling schemes within the the support of  $p(x)$ .

$$C_n^2 \mathbb{E}_{\mathbf{x}_{1:n} \sim p} [\sigma^2(x_n | \{x_i\}_{i=0}^{n-1})] \leq \frac{C_n^2 \gamma_n}{n}$$

[Srinivas et al. \[2012\]](#) derive the bounds on  $\gamma_n$  (see [Table 1](#)) for linear, RBF, and Matèrn kernels on compact and convex sets. Hence, we obtain that for the linear and RBF kernel,  $C_n^2 \gamma_n$  grows sublinearly in  $n$ , i.e.,  $C_n^2 \gamma_n / n \rightarrow 0$  for  $n \rightarrow \infty$ .  $\square$

Kernel	Bounds on $\gamma_n$ for $x \in \mathbb{R}^d$
Linear	$\mathcal{O}(d \log n)$
RBF	$\mathcal{O}((\log n)^{d+1})$
Matèrn $\nu > 1/2$	$\mathcal{O}(n^{\frac{d}{2\nu+d}} \log^{\frac{2\nu}{2\nu+d}}(n))$

Table 1: Bounds on  $\gamma_n$  from [\[Vakili et al., 2021, Theorem 5.\]](#)

**Lemma B.8.** *Let  $p(x)$  denote a distribution with compact support, and assume that  $\mathbf{x}_1, \dots, \mathbf{x}_{n-1}$  are i.i.d. samples from  $p$ . Then, the following holds,*

$$\mathbb{P} \left( \mathbb{E}_{\mathbf{x} \sim p} [C_n \sigma_n(\mathbf{x})] = \mathcal{O} \left( (1 + 1/\sqrt{\delta}) \sqrt{\frac{\gamma_n^{T+1}}{n}} \right), \forall \mathbf{x}_{1:n-1} \right) \geq 1 - \delta, \quad \forall n \in \mathbb{N}.$$

For kernels with a maximum information capacity  $\gamma_n = \mathcal{O}(\text{poly}(\log(n)))$  that grows at most polylogarithmically with  $n$ , we have that

$$\mathbb{P} \left( \mathbb{E}_{\mathbf{x} \sim p} [C_n \sigma_n(x)] \rightarrow 0 \text{ for } n \rightarrow \infty \right) = 1.$$

*Proof.* From [Lemma B.7](#)

$$\mathbb{E}_{\mathbf{x}_{1:n-1} \sim p} \left[ C_n^2 \mathbb{E}_{\mathbf{x} \sim p} [\sigma_n^2(x)] \right] \leq \frac{C_n^2 \gamma_n}{n}. \quad (24)$$

Using the Markov inequality, we get

$$\mathbb{P}(X \geq a) \leq \frac{\mathbb{E}[X]}{a}.$$

Let  $X$  denote  $C_n^2 \mathbb{E}_{\mathbf{x} \sim p} [\sigma_n^2(x)]$ . Then we have for all  $a > 0$ ,

$$\mathbb{P}(C_n^2 \mathbb{E}_{\mathbf{x} \sim p} [\sigma_n^2(x)] \geq a) \leq \frac{\mathbb{E}_{\mathbf{x}_{1:n-1} \sim p} \left[ C_n^2 \mathbb{E}_{\mathbf{x} \sim p} [\sigma_n^2(x)] \right]}{a} \leq \frac{C_n^2 \gamma_n}{na}.$$

Therefore, for  $n \rightarrow \infty$ ,  $C_n^2 \mathbb{E}_{\mathbf{x} \sim p} [\sigma_n^2(x)] \rightarrow 0$  almost surely if  $\frac{C_n^2 \gamma_n}{n} \rightarrow 0$  for  $n \rightarrow \infty$ . Now by definition of  $C_n$  ([Theorem 3.5](#)) and plugging in the choice of  $\beta_n$  ([Lemma 3.6](#)), we have  $\frac{C_n^2 \gamma_n}{n} \propto \frac{\gamma_n^{T+1}}{n}$ . By assumption, we have that  $\gamma_n = \mathcal{O}(\text{poly}(\log(n)))$ , and, thus  $\frac{\gamma_n^{T+1}}{n} = \mathcal{O}(\text{poly}(\log(n))^{T+1}/n) = \mathcal{O}(\text{poly}(\log(n))/n)$ . Hence,  $\frac{C_n^2 \gamma_n}{n} \rightarrow 0$  for  $n \rightarrow \infty$ . For example, for the linear and RBF kernel, we have (see [Table 1](#))

$$\frac{\gamma_n^{T+1}}{n} = \mathcal{O} \left( \frac{d^{T+1} (\log n)^{T+1}}{n} \right) \quad (\text{Linear kernel})$$

$$\frac{\gamma_n^{T+1}}{n} = \mathcal{O} \left( \frac{(\log n)^{(d+1)(T+1)}}{n} \right). \quad (\text{RBF kernel})$$

Now to recover the rate of convergence, let  $v = \mathbb{E}_{x \sim p} [C_n \sigma_n(x)]$ . We study its variance and expectation with respect to  $x_{1:n-1} \sim p$  for a fixed  $n$ . We have

$$\text{Var}[v] \leq \mathbb{E}[v^2] \leq \mathbb{E}_{x_{1:n-1} \sim p} \left[ C_n^2 \mathbb{E}_{x \sim p} [\sigma_n^2(x)] \right] \leq \frac{C_n^2 \gamma_n}{n}.$$

Additionally,  $\mathbb{E}[v] \leq \sqrt{\mathbb{E}[v^2]} \leq C_n \sqrt{\frac{\gamma_n}{n}}$ . Now, we apply the Chebyshev inequality, i.e.,

$$\mathbb{P}(|v - \mathbb{E}[v]| \geq a) \leq \frac{\text{Var}[v]}{a^2} \leq \frac{\mathbb{E}[v^2]}{a^2}.$$

Therefore, for  $a^2 = \frac{\mathbb{E}[v^2]}{\delta}$ , we have with probability at least  $1 - \delta$ ,

$$\begin{aligned} v &\leq \mathbb{E}[v] + a \\ &= \mathbb{E}[v] + \sqrt{\frac{\mathbb{E}[v^2]}{\delta}} \\ &\leq \left(1 + \frac{1}{\sqrt{\delta}}\right) \sqrt{\mathbb{E}[v^2]}. \end{aligned}$$

Next, we plug in the definition of  $v$ , to get

$$\begin{aligned} \mathbb{E}_{x \sim p} [C_n \sigma_n(x)] &\leq \left(1 + \frac{1}{\sqrt{\delta}}\right) \sqrt{\mathbb{E}_{x_{1:n-1} \sim p} \left[ \mathbb{E}_{x \sim p} [C_n^2 \sigma_n^2(x)] \right]} \\ &\leq \left(1 + \frac{1}{\sqrt{\delta}}\right) C_n \sqrt{\frac{\gamma_n}{n}} = \mathcal{O} \left( \left(1 + 1/\sqrt{\delta}\right) \sqrt{\frac{\gamma_n^{T+1}}{n}} \right). \end{aligned}$$

with probability at least  $1 - \delta$ . □

*Proof of Theorem 3.8 (Consistency of HAMBO).* For the GP case we prove that the well calibration assumption, and the Lipschitz continuity of  $f$  and  $\sigma$  are satisfied (see Lemmas 3.6 and 3.7). This allows us to apply Theorem 3.5 and Proposition 3.2, which gives with probability at least  $1 - \delta$  that,

$$J(\pi_e) \geq \tilde{J}(\pi_e) \geq J(\pi_e) - C_n \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \rho^{\pi_e}} [\|\boldsymbol{\sigma}_n(\mathbf{s}, \mathbf{a})\|_2]. \quad (25)$$

To prove consistency, we then only need to show that  $C_n \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \rho^{\pi_e}} [\|\boldsymbol{\sigma}_n(\mathbf{s}, \mathbf{a})\|_2]$  goes to 0 for  $n \rightarrow \infty$ . Since the support of the behavioural policy's state-occupancy measure  $\rho^{\pi_b}$  is compact, and  $\text{supp}(\rho^{\pi_e}) \subseteq \text{supp}(\rho^{\pi_b})$ , we have  $\rho^{\pi_b}(s, a) \geq \hat{C} \rho^{\pi_e}(s, a)$  for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , and some  $\hat{C} > 0$ , i.e., the importance sampling ratio is bounded. We can then write,

$$\begin{aligned} C_n \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \rho^{\pi_e}} [\|\boldsymbol{\sigma}_n(\mathbf{s}, \mathbf{a})\|_2] &\leq \sum_{i=1}^{d_s} \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \rho^{\pi_e}} [C_n \sigma_{n,i}(\mathbf{s}, \mathbf{a})] \\ &= \sum_{i=1}^{d_s} \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \rho^{\pi_b}} \left[ C_n \sigma_{n,i}(\mathbf{s}, \mathbf{a}) \frac{\rho^{\pi_e}(s, a)}{\rho^{\pi_b}(s, a)} \right] \\ &\leq \frac{1}{\hat{C}} \sum_{i=1}^{d_s} \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \rho^{\pi_b}} [C_n \sigma_{n,i}(\mathbf{s}, \mathbf{a})]. \end{aligned}$$

Moreover, by taking a union bound over the dimensions  $1, \dots, d_s$ , Lemma B.8 implies that with probability greater than  $1 - \delta$ , for any set of i.i.d. trajectories,

$$\sum_{i=1}^{d_s} \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \rho^{\pi_e}} [C_n \sigma_{n,i}(\mathbf{s}, \mathbf{a})] = \mathcal{O} \left( d_s \left(1 + \sqrt{d_s/\delta}\right) \sqrt{\frac{\gamma_n^{T+1}}{n}} \right).$$

	Ours		Model-based				Model-free	
	HAMBO-CA	HAMBO-DA1	Rigter et al. [2022]	Yu et al. [2021]	Yu et al. [2020]	Kidambi et al. [2020]	Kumar et al. [2020]	Kostrikov et al. [2022]
HalfCheetah-random	37.1	35.1	39.5	38.8	35.4	25.6	19.6	-
HalfCheetah-medium	66.9	67.9	77.9	54.2	69.5	42.1	49.0	47.4

Table 2: Comparisons on the HalfCheetah from the D4RL benchmark suite. Results of the other algorithms are taken from Rigter et al. [2022].

Consider a sequence  $\{\delta_n\}_{n \geq 0}$  such that  $\lim_{n \rightarrow 0} \delta_n = 0$ , and  $\lim_{n \rightarrow \infty} d_s \left(1 + \sqrt{d_s/\delta_n}\right) \sqrt{\gamma_n^{T+1}/n} = 0$  (e.g.,  $\delta_n = \gamma_n^{-1}$ ), and let  $S_n = \sum_{i=1}^{d_s} C_n \sigma_{n,i}(\mathbf{s}, \mathbf{a})$ . Then we have for all  $\epsilon > 0$

$$\begin{aligned} \sum_{n=0}^{\infty} \mathbb{P}(S_n > \epsilon) &= \sum_{n=0}^{N^*(\epsilon)} \mathbb{P}(S_n > \epsilon) + \sum_{n=N^*(\epsilon)}^{\infty} \mathbb{P}(S_n > \epsilon) \\ &\leq N^*(\epsilon) + \sum_{n=N^*(\epsilon)}^{\infty} \delta_n < \infty, \end{aligned}$$

where,  $N^*(\epsilon)$  is the smallest integer such that  $d_s \left(1 + \sqrt{d_s/\delta}\right) \sqrt{\gamma_n^{T+1}/n} \leq \epsilon$ . This implies that

$$\mathbb{P}\left(\sum_{i=1}^{d_s} \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \rho^{\sigma_\epsilon}} [C_n \sigma_n^i(\mathbf{s}, \mathbf{a})] \rightarrow 0 \text{ for } n \rightarrow \infty\right) = 1.$$

□

## C HAMBO FOR OFFLINE REINFORCEMENT LEARNING

OPE methods are commonly used in offline reinforcement learning (ORL) Levine et al. [2020] to recommend/learn an optimal policy. Moreover, ORL methods also suffer from distribution shifts and are susceptible to overestimation, i.e., overestimating the performance of the recommended policy. Therefore, in principle, a good COPE method can be applied for ORL applications. To this end, we propose a natural modification of HAMBO-CA for ORL.

$$\tilde{J}(\pi^*) := \max_{\pi} \min_{\eta} J_{\tilde{\rho}_\eta}(\pi). \quad (26)$$

Our proposed method induces pessimism with respect to the epistemic uncertainty of the learned transition model to tackle distribution shifts. Similar, to HAMBO-CA, we can also use the HAMBO-DS1 variant to induce pessimism.

We compare our HAMBO-based ORL variants to other ORL algorithms on the OpenAI Gym tasks from the D4RL benchmark Fu et al. [2020]. Specifically, we consider the HalfCheetah environment with data sets generated with a random and a mediocre-performing policy. Our results are presented in table 2.

The max-min optimization in eq (26) is typically very challenging. For our experiments we use the soft actor critic algorithm to train the policy and adversary together (DQN algorithm is used for the HAMBO-DA1 variant).

Note, our proposed ORL algorithms recommend the policy with the best lower bound and not the best expected return (see eq 26). Therefore, in general, they may fail to recommend the optimal policy. This is the price we pay for inducing robustness in our ORL methods. However, in practice (see table 2) we observe that the HAMBO based ORL methods perform competitively to the state of the art in the field.

## References

- Sayak Ray Chowdhury and Aditya Gopalan. On kernelized multi-armed bandits. In *International Conference on Machine Learning*, 2017.
- Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing)*. Wiley-Interscience, 2006.

- 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.
- Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep data-driven reinforcement learning. *arXiv*, 2020.
- 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, 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*, 2018.
- Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel: Model-based offline reinforcement learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q-learning. In *International Conference on Learning Representations (ICLR)*, 2022.
- Volodymyr Kuleshov, Nathan Fenner, and Stefano Ermon. Accurate uncertainties for deep learning using calibrated regression. In *International Conference on Machine Learning (ICML)*, 2018.
- Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- 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.
- Marc Rigter, Bruno Lacerda, and Nick Hawes. Rambo-rl: Robust adversarial model-based offline reinforcement learning. *arXiv*, 2022.
- Niranjan Srinivas, Andreas Krause, Sham M. Kakade, and Matthias W. Seeger. Information-theoretic regret bounds for gaussian process optimization in the bandit setting. *IEEE Transactions on Information Theory*, 2012.
- Sattar Vakili, Kia Khezeli, and Victor Picheny. On information gain and regret bounds in gaussian process bandits. In *International Conference on Artificial Intelligence and Statistics*, 2021.
- Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea Finn, and Tengyu Ma. Mopo: Model-based offline policy optimization. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- Tianhe Yu, Aviral Kumar, Rafael Rafailov, Aravind Rajeswaran, Sergey Levine, and Chelsea Finn. Combo: Conservative offline model-based policy optimization. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.