
Receding Horizon Inverse Reinforcement Learning

Yiqing Xu

Department of Computer Science
National University of Singapore
Singapore, 117417
xuyiqing@comp.nus.edu.sg

Wei Gao

Department of Computer Science
National University of Singapore
Singapore, 117417
qazmichaelgw@gmail.com

David Hsu

Department of Computer Science
National University of Singapore
Singapore, 117417
dyhsu@comp.nus.edu.sg

Abstract

Inverse reinforcement learning (IRL) seeks to infer a cost function that explains the underlying goals and preferences of expert demonstrations. This paper presents *receding horizon inverse reinforcement learning* (RHIRL), a new IRL algorithm for high-dimensional, noisy, continuous systems with black-box dynamic models. RHIRL addresses two key challenges of IRL: scalability and robustness. To handle high-dimensional continuous systems, RHIRL matches the induced optimal trajectories with expert demonstrations *locally* in a receding horizon manner and “stitches” together the local solutions to learn the cost; it thereby avoids the “curse of dimensionality”. This contrasts sharply with earlier algorithms that match with expert demonstrations *globally* over the entire high-dimensional state space. To be robust against imperfect expert demonstrations and system control noise, RHIRL learns a state-dependent cost function “disentangled” from system dynamics under mild conditions. Experiments on benchmark tasks show that RHIRL outperforms several leading IRL algorithms in most instances. We also prove that the cumulative error of RHIRL grows linearly with the task duration.

1 Introduction

Reinforcement learning (RL) has made exciting progress in a range of complex tasks, including real-time game playing [1], visuo-motor control of robots [2], etc. The success, however, often hinges on a carefully crafted cost function [3, 4], which is a major impediment to the wide adoption of RL in practice. Inverse reinforcement learning (IRL) [5] addresses this need by learning a cost function that explains the underlying goals and preferences of expert demonstrations. This work focuses on two key challenges in IRL, *scalability* and *robustness*.

Classic IRL algorithms commonly consist of two nested loops. The inner loop approximates the optimal control policy for a hypothesized cost function, while the outer loop updates the cost function by comparing the induced policy with expert demonstrations. The inner loop solves the (forward) reinforcement learning or optimal control problem, which is in itself a challenge for complex high-dimensional systems. Many interesting ideas have been proposed for IRL, including, e.g., maximum entropy learning [6, 7], guided cost learning [8], and adversarial learning [9]. See Figure 1 for illustrations. They all try to match a *globally* optimal approximate policy with expert demonstrations over the entire system state space or a sampled approximation of it. This is impractical for high-dimensional continuous systems and is a fundamental impediment to scalability. Like RL, IRL also

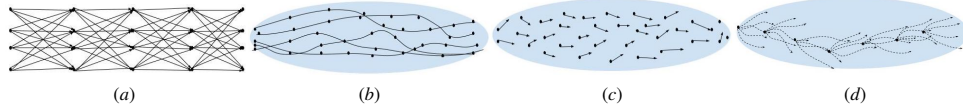


Figure 1: A comparison of RHIRL and selected IRL algorithms. They all try to match the policy induced by the learned cost with expert demonstrations. (a) MaxEnt matches the exact feature count over the entire system state space. (b) REIRL, GCL, and GAN-GCL match the approximate feature count over sampled state trajectories globally over the entire task duration. (c) GAIL, AIRL, ... discriminate the sampled state or state-action distributions. (d) RHIRL matches control sequences *locally* along demonstrated trajectories in a receding horizon manner.

suffers from the “curse of dimensionality”. To scale up, RHIRL computes locally optimal policies with receding horizons rather than a globally optimal policy and then matches them with expert demonstrations *locally* in succession (Figure 1d). The local approximation and matching substantially mitigate the impact of high-dimensional space and improve the sample efficiency of RHIRL, at the cost of a local rather than a global solution. So RHIRL trades off optimality for scalability and provides a novel alternative to the existing approaches.

Another important concern of IRL is noise in expert demonstrations and system control. Human experts may be imperfect for various reasons and provide good, but still suboptimal demonstrations. Further, the system may fail to execute the commanded actions accurately because of control noise. We want to learn a cost function that captures the expert’s underlying intentions rather than the imperfectly executed actions. We recognize that there are a number of IRL works on learning from sub-optimal/fail demonstrations [10, 11, 12]. There is a subtle, but key difference between (i) demonstrations intended to be suboptimal and (ii) optimal demonstrations corrupted by noise. The existing works mentioned above address (i). RHIRL addresses (ii). To learn the true intentions from the noise corrupted demonstrations, RHIRL relies on a simplifying assumption: the cost function is linearly separable with two components, one state-dependent and one control-dependent. Many interesting systems in practice satisfy the assumption, at least, approximately [13]. RHIRL then learns the state-dependent cost component, which is disentangled from the system dynamics [9] and agnostic to noise.

2 Related Work

IRL can be viewed as an indirect approach to imitation learning. It learns a cost function, which induces an optimal policy that matches with expert demonstrations. In contrast, behavior cloning (BC) is a direct approach. It assumes independence among all demonstrated state-action pairs and learns a policy that maps states to actions through supervised learning of state-action pairs from expert demonstrations. The simplicity of BC is appealing. However, it usually requires large amounts of data to learn well and suffers from covariate shift [14]. IRL is more data-efficient. Further, it produces a cost function, which explains the expert demonstrations and potentially transfers to other systems with different dynamics. These benefits, however, come at the expense of greater computational complexity.

Classic IRL algorithms learn a cost function iteratively in a double-loop: the outer loop updates a hypothesized cost function, and the inner loop solves the forward RL problem for an optimal policy and matches it with expert demonstrations. Various methods have been proposed [6, 8, 7, 15], but they all seek a globally optimal solution over the entire state space (Figure 1a) [6, 7] or the entire task duration (Figure 1b) [8, 15]. As a result, they face significant challenges in scalability and must make simplifications, such as locally linear dynamics [8]. Recent methods use the generative adversarial network (GAN) [16] to learn a discriminator that differentiates between the state or state-action distribution induced by the learned cost and that from expert demonstrations [9, 17, 18, 19, 20, 21, 22, 23, 24]. We view this global matching as a major obstacle to scalability. In addition, GAN training is challenging and faces difficulties with convergence.

RHIRL stands in between BC and existing IRL algorithms as a novel midpoint by trading off optimality for scalability. BC performs local matching at each demonstrated state-action pair, treating all of them independently. Existing IRL algorithms perform global matching over the entire state space or sampled trajectories from it. RHIRL follows the standard IRL setup. To tackle the challenge of scalability for high-dimensional continuous systems, RHIRL borrows ideas from receding horizon

control [25]. It solves for locally optimal control sequences with receding horizons and learns the cost function by “stitching” together a series of locally optimal solutions to match the global state distribution of expert demonstrations (Figure 1d).

Another challenge to IRL is suboptimal expert demonstrations and system control noise. Several methods learn an auxiliary score or ranking to reweigh the demonstrations, in order to approximate the underlying optimal expert distribution [26, 27, 22]. RHIRL does not attempt to reconstruct the optimal expert demonstrations. It explicitly models noise in control actions and matches the noisy control with expert demonstrations, in order to learn from the intended, rather than the executed expert actions.

3 Receding Horizon Inverse Reinforcement Learning

3.1 Overview

Consider a continuous dynamical system with noisy control:

$$x_{t+1} = f(x_t, v_t), \quad (1)$$

where $x_t \in \mathcal{R}^n$ is the state, $v_t \in \mathcal{R}^m$ is the control at time t , and the initial system state at $t = 0$ follows a distribution μ . To account for noise in expert demonstrations and in system control, we assume that v_t is a random variable following the Gaussian distribution $\mathcal{N}(v_t|u_t, \Sigma)$, with mean u_t and covariance Σ . We can control u_t directly, but not v_t , because of noise. The state-transition function f captures the system dynamics. RHIRL represents f as a black-box simulator and does not require its analytic form. Thus, we can accommodate arbitrary complex nonlinear dynamics. To simplify the presentation, we assume that the system dynamics is deterministic. We sketch the extension to stochastic dynamics at the end of the section, the full proof is given in Appendix C.

In RL, we are given a cost function and want to find a control policy that minimizes the expected total cost over time under the dynamical system. In IRL, we are not given the cost function, but instead, a set \mathcal{D} of expert demonstrations. Each demonstration is a trajectory of states visited by the expert over time: $\tau = (x_0, x_1, x_2, \dots, x_{T-1})$ for a duration of T steps.

We assume that the expert chooses the actions to minimize an unknown cost function and want to re this cost from the demonstrations. Formally, suppose that the cost function is parameterized by θ . RHIRL aims to learn a cost function that minimizes the loss $\mathcal{L}(\theta; \mathcal{D})$, which measures the difference between the demonstration trajectories and the optimal control policy induced by the cost function with parameters θ .

RHIRL performs this minimization iteratively, using the gradient $\partial\mathcal{L}/\partial\theta$ to update θ . In iteration t , let x_t be the system state at time t and \mathcal{D}_t be the set of expert sub-trajectories starting at time t and having a duration of maximum K steps. We use the current cost to perform receding horizon control (RHC) at x_t , with time horizon K , and then update the cost by comparing the resulting state trajectory distribution with the demonstrations in \mathcal{D}_t . See Algorithm 1 for a sketch.

First, sample M control sequences, each of length K (line 6). We assume that the covariance Σ is known. If it is unknown, we set it to be identity by default. For our experiment results reported in Table 1 and Table 2, Σ is unknown and is approximated by a constant factor of the identity matrix (grid search is performed to determine this constant factor). We show through experiments that the learned cost function is robust over different noise settings (Section 4.3).

Next, we apply model-predictive path integral (MPPI) control [28] at x_t . MPPI provides an analytic solution for the optimal control sequence distribution and the associated state sequence distribution, which allow us to estimate the gradient $\partial\mathcal{L}/\partial\theta$ efficiently through importance sampling (lines 7–8) and update cost function parameters θ (line 9). Finally, we execute the computed optimal control (line 10–12) and update the mean control input for the next iteration (line 13). We would like to emphasize that we only use the simulator to sample rollouts and evaluate the current cost function. Fixed expert demonstrations \mathcal{D} are given upfront. Unlike DAgger, we do not query the expert online during training.

Our challenge is to uncover a cost function that captures the expert’s intended controls $(u_0, u_1, \dots, u_{T-1})$, even though they were not directly observed, because of noise, and do so in a scalable and robust manner for high-dimensional, noisy systems.

Algorithm 1 RHIRL

```
1: Initialize  $\theta$  randomly.
2: for  $i = 1, 2, 3, \dots$  do
3:   Sample  $x_0$  from  $\mu$ .
4:   Initialize control sequence  $U$  of length  $K$  to  $(0, 0, \dots)$ .
5:   for  $t = 0, 1, 2, \dots, T - 1$  do
6:     Sample  $M$  control sequences  $V_j$  for  $j = 1, 2, \dots, M$ , according to  $\mathcal{N}(V|U, \Sigma)$ .
7:     Compute the importance weight  $w_j$  for each  $V_j$ , using the state cost  $S(V_j, x_t; \theta)$ . See equation (11).
8:     Compute  $\frac{\partial}{\partial \theta} \mathcal{L}(\theta; \mathcal{D}_t, x_t)$  using  $\mathcal{D}_t$  and  $V_j$  with weight  $w_j$ , for  $j = 1, 2, \dots, M$ . See equation (12).
9:      $\theta \leftarrow \theta - \alpha \frac{\partial \mathcal{L}}{\partial \theta}$ .
10:     $U \leftarrow \sum_{j=1}^M w_j V_j$ .
11:    Sample  $v_t$  from  $\mathcal{N}(v|u, \Sigma)$ , where  $u$  is the first element in the control sequence  $U$ .
12:     $x_{t+1} \leftarrow f(x_t, v_t)$ .
13:    Remove  $u$  from  $U$ . Append 0 at the end of  $U$ .
14:  end for
15: end for
```

We develop three ideas: structuring the cost function, matching locally with expert demonstrations, and efficient computation of the gradient $\partial \mathcal{L} / \partial \theta$, which are described next.

3.2 Robust Cost Functions

To learn a cost function robust against noise, we make a simplifying assumption that linearly separates the one-step cost into two components: a state cost $g(x; \theta)$ parameterized by θ and a quadratic control cost $u^T \Sigma^{-1} u$. Despite the simplification, this cost function models a wide variety of interesting systems in practice [13]. It allows RHIRL to learn a state cost $g(x; \theta)$, independent of control noise (Section 3.4), and thus generalize over different noise distributions (Section 4.3).

Suppose that $V = (v_0, v_1, \dots, v_{K-1})$ is a control sequence of length K , conditioned on the input $U = (u_0, u_1, \dots, u_{K-1})$. We apply V to the dynamical system in (1) with start state x_0 and obtain a state sequence $\tau = (x_0, x_1, \dots, x_K)$ with $x_k \sim p(x_k | x_{k-1}, v_{k-1})$ for $k = 1, 2, \dots, K$. Define the total cost of V as

$$J(V, x_0; \theta) = \sum_{k=0}^K g(x_k; \theta) + \sum_{k=0}^{K-1} \frac{\lambda}{2} u_k^T \Sigma^{-1} u_k, \quad (2)$$

where $\lambda \geq 0$ is a constant weighting the relative importance between the state and control costs. For convenience, define also the total state cost of V as

$$S(V, x_0; \theta) = \sum_{k=0}^K g(x_k; \theta). \quad (3)$$

While S is defined in terms of the control sequence V , it only depends on the corresponding state trajectory τ . This is very useful, as the training data contains state and not control sequences explicitly.

3.3 Local Control Sequence Matching

To minimize the loss \mathcal{L} , each iteration of RHIRL applies RHC with time horizon K under the current cost parameters θ and computes locally optimal control sequences of length K . In contrast, classic IRL algorithms, such as MaxEnt [6], perform global optimization over the entire task duration T in the inner loop. While RHC sacrifices global optimality, it is much more scalable and enables RHIRL to handle high-dimensional continuous systems. We use the hyperparameter K to trade off optimality and scalability.

Specifically, we use MPPI [28] to solve for an optimal control sequence distribution at the current start state x_t in iteration t . The main result of MPPI suggests that the optimal control sequence distribution Q minimizes the “free energy” of the dynamical system and this free energy can be calculated from the cost of the state trajectory under Q . Mathematically, the probability density $p_Q(V^*)$ can be expressed as a function of the state cost $S(V, x_t; \theta)$, with respect to a Gaussian “base” distribution

$B(U_B, \Sigma)$ that depends on the control cost:

$$p_Q(V^*|U_B, \Sigma, x_t; \theta) = \frac{1}{Z} p_B(V^*|U_B, \Sigma) \exp\left(-\frac{1}{\lambda} S(V^*, x_t; \theta)\right), \quad (4)$$

where Z is the partition function. For the quadratic control cost in (37), we have $U_B = (0, 0, \dots)$ [28]. Intuitively, the expression in (4) says that the probability density of V^* is the product of two factors, one penalizing high control cost and one penalizing high state cost. So controls with large values or resulting in high-cost states occur with probability exponentially small.

Equation (4) provides the optimal control sequence distribution under the current cost. Suppose that the control sequences for expert demonstrations \mathcal{D}_t follow a distribution E . We define the loss $\mathcal{L}(\theta; \mathcal{D}_t, x_t)$ as the KL-divergence between the two distributions:

$$\mathcal{L}(\theta; \mathcal{D}_t, x_t) = D_{\text{KL}}(p_E(V|x_t) \parallel p_Q(V|U_B, \Sigma, x_t; \theta)), \quad (5)$$

which RHIRL seeks to minimize in each iteration. While the loss \mathcal{L} is defined in terms of control sequence distributions, the expert demonstrations \mathcal{D} provide state information and not control information. However, each control sequence V induces a corresponding state sequence τ for a given start state x_0 , and τ determines the cost of V according to (37). We show in the next subsection that $\partial \mathcal{L} / \partial \theta$ can be computed efficiently using only state information from \mathcal{D} . This makes RHIRL appealing for learning tasks in which control labels are difficult or impossible to acquire.

3.4 Gradient Optimization

To simplify notations, we remove the explicit dependency on x_t and \mathcal{D}_t and assume that in iteration t of RHIRL, all control sequences are applied with x_t as the start state and expert demonstrations come from \mathcal{D}_t . We have

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \theta} &= \frac{\partial}{\partial \theta} \int p_E(V) \log \frac{p_E(V)}{p_Q(V|U_B, \Sigma; \theta)} dV \\ &= \int p_E(V) \left(\frac{1}{\lambda} \frac{\partial}{\partial \theta} S(V; \theta) \right) dV - \int p_Q(V|U_B, \Sigma; \theta) \left(\frac{1}{\lambda} \frac{\partial}{\partial \theta} S(V; \theta) \right) dV \end{aligned} \quad (6)$$

The first line follows directly from the definition, and the derivation for the second line is available in Appendix A.

We estimate the two integrals in (6) through sampling. For the first integral, we can use the expert demonstrations as samples. For the second integral, we cannot sample p_Q directly, as the optimal control distribution Q is unknown in advance. Instead, we sample from a known Gaussian distribution with density $p(V|U, \Sigma)$ and apply importance sampling so that

$$\mathbb{E}_{p_Q(V|U_B, \Sigma; \theta)}[V] = \mathbb{E}_{p(V|U, \Sigma)}[w(V)V]. \quad (7)$$

The importance weight $w(V)$ is given by

$$w(V) = \frac{p_Q(V|U_B, \Sigma; \theta)}{p(V|U, \Sigma)} = \frac{p_Q(V|U_B, \Sigma; \theta)}{p_B(V|U_B, \Sigma)} \times \frac{p_B(V|U_B, \Sigma)}{p(V|U, \Sigma)} \quad (8)$$

To simplify the first ratio in (8), we substitute in the expression for p_Q from (4):

$$\frac{p_Q(V|U_B, \Sigma; \theta)}{p_B(V|U_B, \Sigma)} = \frac{1}{Z} \exp\left(-\frac{1}{\lambda} S(V; \theta)\right) \quad (9)$$

We then simplify the second ratio, as both are Gaussian distributions with the same covariance matrix Σ :

$$\frac{p_B(V|U_B, \Sigma)}{p(V|U, \Sigma)} \propto \exp\left(-\sum_{k=0}^{K-1} (u_k - u_k^B)^\top \Sigma^{-1} v_k\right), \quad (10)$$

where u_k and v_k , $k = 0, 1, \dots, K-1$ are the mean controls and the sampled controls from $p(V|U, \Sigma)$, respectively. Similarly, u_k^B , $k = 0, 1, \dots, K-1$ are the mean controls for the base distribution, and they are all 0 in our case. Substituting (9) and (10) into (8), we have

$$w(V) \propto \exp\left(-\frac{1}{\lambda} \left(S(V; \theta) + \lambda \sum_{k=0}^{K-1} u_k^\top \Sigma^{-1} v_k \right) \right) \quad (11)$$

For each sampled control sequence V , we evaluate the expression in (11) and normalize over all samples to obtain $w(V)$.

To summarize, we estimate $\partial\mathcal{L}/\partial\theta$ through sampling:

$$\frac{\partial}{\partial\theta}\mathcal{L}(\theta; \mathcal{D}_t, x_t) \approx \frac{1}{N} \sum_{i=1}^N \frac{1}{\lambda} \frac{\partial}{\partial\theta} S(V_i, x_t; \theta) - \frac{1}{M} \sum_{j=1}^M \frac{1}{\lambda} w(V_j) \frac{\partial}{\partial\theta} S(V_j, x_t; \theta), \quad (12)$$

where $V_i, i = 1, \dots, N$ are the control sequences for the expert demonstrations in \mathcal{D}_t and $V_j, j = 1, 2, \dots, M$ are the sampled control sequences. Equation (12) connects $\partial\mathcal{L}/\partial\theta$ with $\partial S/\partial\theta$. The state cost function S is represented as a shallow neural network, and its derivative can be obtained easily through back-propagation. To evaluate $\frac{\partial}{\partial\theta} S(V_i, x_t; \theta)$, we do not actually use the expert control sequences, as they are unknown. We use the corresponding state trajectories in \mathcal{D}_t directly, as the state cost depends only on the visited states. See equation (3).

Finally, we approximate the optimal mean control sequence according to (7):

$$U = \mathbb{E}_{p_Q(V|U_B, \Sigma, x_t; \theta)}[V] \approx \sum_{j=1}^M w(V_j) V_j. \quad (13)$$

The first element in the control sequence U is the chosen control for the current time step t . We remove the first element from U and append zero at the end. The new control sequence is then used as the mean for the sampling distribution in the next iteration.

3.5 Analysis

Since RHIRL performs local optimization sequentially over many steps, one main concern is error accumulation over time. For example, standard behavior cloning has one-step error ϵ and cumulative error $O(T^2\epsilon)$ over T steps, because of covariate shift [14]. The DAgger algorithm reduces the error to $O(T\epsilon)$ by querying the expert repeatedly during online learning [14]. We prove a similar result for RHIRL, which uses offline expert demonstrations only. In iteration t of RHIRL, let $p_E(V_t|x_t)$ be the K -step expert demonstration distribution and $p_Q(V_t|U_B, \Sigma, x_t; \theta)$ be the computed K -step optimal control distribution for some fixed control parameters θ . RHIRL minimizes the KL-divergence between these two distributions in each iteration. Let $p_E(x)$ be the state marginal distribution of expert demonstrations and $p_{\text{RHC}}(x; \theta)$ be the state marginal distribution of the computed RHC policy under θ over the entire task duration T . Intuitively, we want the control policy under the learned cost to visit states similar to those of expert demonstrations. In other words, $p_{\text{RHC}}(x; \theta)$ and $p_E(x)$ are close.

Theorem 3.1. *If $D_{\text{KL}}(p_E(V_t|x_t) \parallel p_Q(V_t|U_B, \Sigma, x_t; \theta)) < \epsilon$ for all $t = 0, 1, \dots, T-1$, then $D_{\text{TV}}(p_E(x), p_{\text{RHC}}(x; \theta)) < T\sqrt{\epsilon/2}$.*

The theorem says that RHIRL’s cumulative error, measured in total variation distance D_{TV} between $p_E(x)$ and $p_{\text{RHC}}(x; \theta)$, grows linearly with T . The proof consists of three steps. First, in each iteration t , if the KL-divergence between two control sequence distributions are bounded by ϵ , so is the KL-divergence between the two corresponding state distributions induced by control sequences. Next, we show that the KL-divergence between the state distributions over two successive time steps are bounded by the same ϵ . Finally, we switch from KL-divergence to total variation distance and apply the triangle inequality to obtain the final result. The full proof is given in Appendix B. Since RHC performs local optimization in each iteration, we cannot guarantee global optimality. However, the theorem indicates that unlike standard behavior cloning, the cumulative error of RHIRL grows linearly and not quadratically in T . This shows one advantage of IRL over behavior cloning from the theoretical angle.

Given a control policy V with the resulting state marginal distribution $p_V(x)$, we can calculate the expected total cost of V by integrating the one-step cost over p_V . Now suppose that the one-step cost is bounded. Theorem 3.1 then implies that the regret in total cost, compared with the expert policy, also grows linearly in T .

3.6 Extension to Stochastic Dynamics

Suppose that the system dynamics is stochastic: $x_{t+1} = f(x_t, v_t, \omega_t)$, where ω_t is a random variable that models the independent system noise. RHIRL still applies, with modifications. We redefine






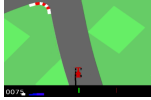
| | | | | | |
|---|---|---|---|---|---|
| Pendulum-v0 | LunarLander-v2 | Hopper-v2 | Walker2d-v2 | Ant-v2 | CarRacing-v0 |
|  |  |  |  |  |  |
| $ S $: 3-D | 8-D | 11-D | 17-D | 111-D | 96×96 (image) |
| $ A $: 1-D | 2-D | 3-D | 6-D | 8-D | 3-D |
| T : 100 | 250 | 1000 | 1000 | 1000 | 1000 |

Figure 2: Benchmark tasks. $|S|$ and $|A|$ denote the dimensions of the state space and the action space, respectively. T denotes the task horizon.

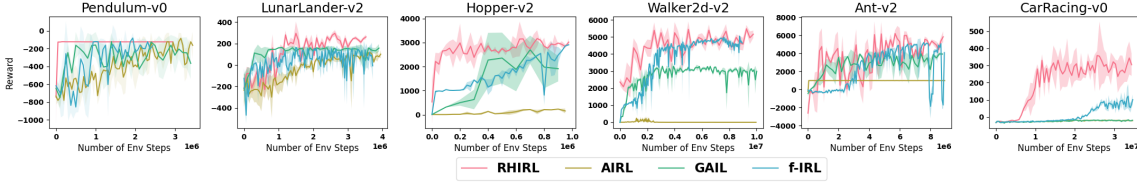


Figure 3: Learning curves for RHIRL and other methods.

the total cost functions $\tilde{J}(V, x_0; \theta)$ and $\tilde{S}(V, x_0; \theta)$ by taking expectation over the system noise distribution. When calculating the importance weight $\tilde{w}(V)$ in (11), we sample over the noise distribution to estimate the expected total state cost. Finally, we may need more samples when estimating the gradient in (12), because of the increased variance due to the system noise. The full proof is given in Appendix C. We acknowledge that the experiments in this work are of deterministic dynamics, we leave the empirical results on RHIRL’s extension to stochastic dynamics to future work.

4 Experiments

We investigate two main questions. Does RHIRL scale up to high-dimensional continuous control tasks? Does RHIRL learn a robust cost function under noise?

4.1 Setup

We compare RHIRL with two leading IRL algorithms, namely AIRL [9] and f -IRL [20], and one imitation learning algorithm, GAIL [17]. In particular, f -IRL is a recent algorithm that achieves leading performance on high-dimensional control tasks. We use the implementation of AIRL, GAIL, and f -IRL from the f -IRL’s official repository along with the reported hyperparameters [20], whenever possible. We also perform hyperparameter search on a grid to optimize the performance of every method on every task. The specific hyperparameter settings used are reported in Appendix D.2.

Our benchmark set consists of six continuous control tasks (Figure 2) from OpenAI Gym [29], with increasing sizes of state and action spaces. For the most complex task, CarRacing, the input consists of 96×96 raw images, resulting in an enormous state space that poses great challenges [30]. To our knowledge, RHIRL is the first IRL algorithm to attempt such a high-dimensional space. For fair comparison, we customize all tasks to the fixed task horizon settings (Figure 2. See Appendix D.2 for details on task parameter settings.

We use WorldModel [30] to generate expert demonstration data for CarRacing and use SAC [31] for the other tasks. We add Gaussian noise to the input controls and collect expert demonstrations at different control noise levels. The covariance of the control noise is *unknown* to all methods, including, in particular, RHIRL.

To measure the performance of the learned cost function and policy, we score its induced optimal policy using the ground-truth cost function. For ease of comparison with the literature, we use negated cost values, i.e., rewards, in all reported results. Higher values indicate better performance. Each experiment is repeated 10 times to estimate the performance variance.

4.2 Scalability

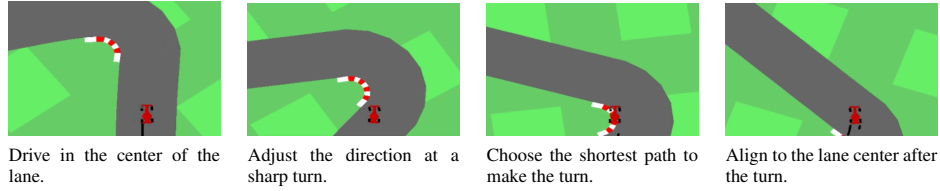


Figure 4: Driving behaviors learned by RHIRL for CarRacing-v0, with pixel-level raw image input.

We compare RHIRL with other methods, first in noise-free environments (Figure 3) and then with increasing noise levels (Table 1).

Figure 3 shows the learning curve of each methods in noise-free environments. Overall, RHIRL converges faster and achieves higher return, especially for tasks with higher state space dimensions. This improved performance suggests that the benefit of local optimization adopted by RHIRL outweighs its potential limitations.

Table 1 shows the final performance of all methods at various noise levels. RHIRL clearly outperforms AIRL and GAIL in all experiments. So we focus our discussion on comparison with f -IRL. In noise-free environments, RHIRL and f -IRL perform comparably on most tasks. On CarRacing, the most challenging task, RHIRL performs much better. RHIRL manages to learn the critical driving behaviors illustrated in Figure 4, despite the high-dimensional image input. However, RHIRL does not manage to learn to drive fast enough. That is the main reason why it under-performs the expert. In comparison, f -IRL only learns to follow a straight lane after a large number of environment steps, and still fails to make a sharp turn after 3.0×10^7 environment steps. In the noisy environments, the advantage of RHIRL over f -IRL is more pronounced even on some of the low-dimensional tasks, because RHIRL accounts for the control noise explicitly in the cost function.

Next we evaluate RHIRL and other methods for robustness under noise. A robust cost function encodes the expert’s true intent. It is expected to yield consistent performance over different noise levels, regardless of noise in expert demonstrations. For each task, a cost function is learned in a noise-free environment and is then used to re-optimize a policy in the corresponding noisy environments. Specifically for GAIL, since it learns a policy and does not recover the associated cost function, we directly apply the learned policy in noisy environments.

Table 1: Performance comparison of RHIRL and other methods. The performance is reported as the ratio of the learned policy’s average return and the expert’s average return. The absolute average returns and the standard deviations are reported in Appendix D.4. Negative ratios are clipped to 0. The two numbers under the name of each environment indicate the dimensions of the state space and the action space, respectively.

| | | No Noise $\Sigma = 0$ | Mild Noise $\Sigma = 0.2$ | High Noise $\Sigma = 0.5$ |
|--------------------------------|----------|--------------------------|------------------------------|------------------------------|
| Pendulum 3, 1 | Expert | -154.69 \pm 67.61 | -156.50 \pm 70.72 | -168.54 \pm 80.89 |
| | RHIRL | 1.06 | 1.07 | 1.08 |
| | f -IRL | 1.07 | 1.06 | 0.93 |
| | AIRL | 1.05 | 0.94 | 0.91 |
| | GAIL | 0.88 | 0.89 | 0.80 |
| Lunarlander 8, 2 | Expert | 232.00 \pm 86.12 | 222.65 \pm 56.35 | 164.52 \pm 16.79 |
| | RHIRL | 1.05 | 1.04 | 1.20 |
| | f -IRL | 0.76 | 0.63 | 0.74 |
| | AIRL | 0.74 | 0.60 | 0.58 |
| | GAIL | 0.72 | 0.56 | 0.60 |
| Hopper 11, 3 | Expert | 3222.48 \pm 390.65 | 3159.72 \pm 520.00 | 2887.72 \pm 483.93 |
| | RHIRL | 0.95 | 0.98 | 0.96 |
| | f -IRL | 0.96 | 0.82 | 0.43 |
| | AIRL | 0.01 | 0.01 | 0.01 |
| | GAIL | 0.82 | 0.50 | 0.24 |
| Walker2d 17, 6 | Expert | 4999.47 \pm 55.99 | 4500.43 \pm 114.48 | 3624.48 \pm 95.05 |
| | RHIRL | 0.99 | 0.99 | 0.95 |
| | f -IRL | 0.99 | 0.82 | 0.78 |
| | AIRL | 0.00 | 0.00 | 0.00 |
| | GAIL | 0.50 | 0.64 | 0.51 |
| Ant 111, 8 | Expert | 5759.22 \pm 173.57 | 3257.37 \pm 501.95 | 252.62 \pm 91.44 |
| | RHIRL | 0.86 | 0.93 | 0.91 |
| | f -IRL | 0.87 | 0.80 | 0.78 |
| | AIRL | 0.17 | 0.33 | 0.00 |
| | GAIL | 0.48 | 0.40 | 0.00 |
| CarRacing 96 \times 96, 3 | Expert | 903.25 \pm 0.23 | 702.01 \pm 0.3 | 281.12 \pm 0.34 |
| | RHIRL | 0.40 | 0.29 | 0.19 |
| | f -IRL | 0.09 | 0.03 | 0.00 |
| | AIRL | 0.00 | 0.00 | 0.00 |
| | GAIL | 0.00 | 0.01 | 0.00 |

4.3 Robustness

Table 2 shows that noise causes performance degradation in all methods. However, RHIRL is more robust in comparison.

For simple tasks, Pendulum and Lunarlander, RHIRL and f -IRL perform consistently well across different noise levels, while GAIL and AIRL fail to maintain their good performance, when the noise level increases. For the more challenging tasks, Hopper and Walker, RHIRL’s performance degrades mildly, and f -IRL suffers more significant performance degradation.

It is worth noting that the expert demonstrations used to train the transferred cost function are from the perfect system. Therefore, some expert actions and states may no longer be optimal or feasible in a highly noisy environment. Moreover, the cost function trained in the perfect system cannot reason about the long-term consequences of an action in a high noise environment. Therefore, it is challenging for the learned cost function to be robust to a highly noisy environment, as capturing the true intention of the expert is difficult.

4.4 Effect of Receding Horizon K

The receding horizon K allows RHIRL to trade off optimality and efficiency. To study its effect, we use Pendulum with $K \in \{5, 20, 30, 75\}$. When K is small, RHIRL learns fast with fewer environment steps, but may be unstable. For $K = 5$, the performance jumps after the first few steps, but fluctuates drastically, as K is much smaller than the number of look-ahead steps for the pendulum to swing up. The learned cost function tends to overfit to local behaviors and forget the useful reward signals learned previously. With larger K values (20 or 30), the performance of the cost function learned is more stable, without much loss in efficiency. However, when K is too large, e.g., $K = 75$, RHIRL becomes less efficient and takes longer to converge, because that excessive K values enlarge the variance of sampled control sequences.

5 Conclusion

RHIRL is a scalable and robust IRL algorithm for high-dimensional, noisy, continuous systems. Our experiments show that RHIRL outperforms several leading IRL algorithms on multiple benchmark tasks, especially when expert demonstrations are noisy. RHIRL’s choice of local rather than global optimization is an important issue that deserves further investigation. Overall, we view this as an interesting trade-off between scalability and optimality. While this trade-off is well known in reinforcement learning, optimal control, and general optimization problems, it is mostly unexplored in IRL. Further, local optimization may tie the learned cost with the optimizer. It would be interesting to examine whether the learned cost transfers easily to other domains with different optimizers. We are keen to investigate these important issues and their implications to IRL as our next step.

Table 2: Robustness of RHIRL and other methods under noise. The cost functions are learned in noise-free environments and evaluated in noisy environments. The performance is measured as the ratio between the average return of an re-optimized policy in a noisy environment and the expert’s average return in the corresponding noise-free environment. The absolute average returns and the standard deviations are reported in Appendix D.4. The negative ratios are clipped to 0.

| | | Noise Level Σ | | |
|--------------------------------|----------|----------------------|-------------|-------------|
| | | 0.0 | 0.2 | 0.5 |
| Pendulum 3, 1 | Expert | -154.69 \pm 67.61 | — | — |
| | RHIRL | 1.06 | 1.07 | 1.06 |
| | f -IRL | 1.08 | 0.90 | 0.85 |
| | AIRL | 1.05 | 0.79 | 0.67 |
| | GAIL | 0.88 | 0.71 | 0.62 |
| LunarLander 8, 2 | Expert | 232.00 \pm 86.12 | — | — |
| | RHIRL | 1.05 | 0.89 | 0.76 |
| | f -IRL | 0.76 | 0.53 | 0.44 |
| | AIRL | 0.74 | 0.14 | 0.10 |
| | GAIL | 0.72 | 0.44 | 0.34 |
| Hopper 11, 3 | Expert | 3222.48 \pm 390.65 | — | — |
| | RHIRL | 0.95 | 0.80 | 0.67 |
| | f -IRL | 0.96 | 0.65 | 0.62 |
| | AIRL | 0.01 | 0.01 | 0.00 |
| | GAIL | 0.82 | 0.07 | 0.06 |
| Walker 17, 6 | Expert | 4999.47 \pm 55.99 | — | — |
| | RHIRL | 0.99 | 0.80 | 0.69 |
| | f -IRL | 0.99 | 0.60 | 0.22 |
| | AIRL | 0.00 | 0.28 | 0.36 |
| | GAIL | 0.50 | 0.02 | 0.02 |
| Ant 111, 8 | Expert | 5759.22 \pm 173.57 | — | — |
| | RHIRL | 0.86 | 0.55 | 0.15 |
| | f -IRL | 0.87 | 0.35 | 0.08 |
| | AIRL | 0.17 | 0.15 | 0.00 |
| | GAIL | 0.48 | 0.00 | 0.00 |
| CarRacing 96 \times 96, 3 | Expert | 903.25 \pm 0.23 | — | — |
| | RHIRL | 0.40 | 0.29 | 0.12 |
| | f -IRL | 0.09 | 0.02 | 0.00 |
| | AIRL | 0.00 | — | — |
| | GAIL | 0.00 | — | — |

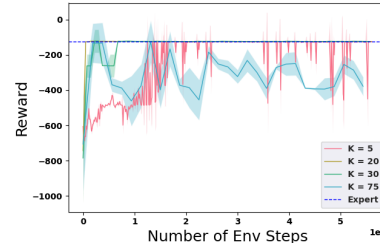


Figure 5: The effect of receding horizon K on the performance of RHIRL on the Pendulum task.

References

- [1] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin A. Riedmiller. Playing atari with deep reinforcement learning. *CoRR*, 2013.
- [2] Ulrich Viereck, Andreas Pas, Kate Saenko, and Robert Platt. Learning a visuomotor controller for real world robotic grasping using simulated depth images. 2017.
- [3] A.Y. Ng, D. Harada, and S. Russell. Policy invariance under reward transformations: Theory and application to reward shaping. 1999.
- [4] W Bradley Knox, Alessandro Allievi, Holger Banzhaf, Felix Schmitt, and Peter Stone. Reward (mis) design for autonomous driving. arXiv:2104.13906, 2021.
- [5] Andrew Y. Ng and Stuart J. Russell. Algorithms for inverse reinforcement learning. 2000.
- [6] Brian D. Ziebart, Andrew Maas, J. Andrew Bagnell, and Anind K. Dey. Maximum entropy inverse reinforcement learning. 2008.
- [7] Markus Wulfmeier, Peter Ondruska, and Ingmar Posner. Maximum entropy deep inverse reinforcement learning, 2016.
- [8] Chelsea Finn, Sergey Levine, and Pieter Abbeel. Guided cost learning: Deep inverse optimal control via policy optimization. 2016.
- [9] Justin Fu, Katie Luo, and Sergey Levine. Learning robust rewards with adversarial inverse reinforcement learning. 2018.
- [10] Kyriacos Shiarlis, João V. Messias, and Shimon Whiteson. Inverse reinforcement learning from failure. In *AAMAS*, 2016.
- [11] Yueh-Hua Wu, Nontawat Charoenphakdee, Han Bao, Voot Tangkaratt, and Masashi Sugiyama. Imitation learning from imperfect demonstration. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *ICML*, volume 97 of *Proceedings of Machine Learning Research*, pages 6818–6827. PMLR, 09–15 Jun 2019.
- [12] Daniel S. Brown, Wonjoon Goo, Prabhat Nagarajan, and Scott Niekum. Extrapolating beyond suboptimal demonstrations via inverse reinforcement learning from observations, 2019.
- [13] W.B. Powell. *Approximate Dynamic Programming: Solving the Curses of Dimensionality*. Wiley-Interscience, 2007.
- [14] Stephane Ross, Geoffrey J. Gordon, and J. Andrew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. 2011.
- [15] Abdeslam Boularias, Jens Kober, and Jan Peters. Relative entropy inverse reinforcement learning. 2011.
- [16] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. 2014.
- [17] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. 2016.
- [18] Ahmed H. Qureshi, Byron Boots, and Michael C. Yip. Adversarial imitation via variational inverse reinforcement learning. 2019.
- [19] Seyed Kamyar Seyed Ghasemipour, Richard Zemel, and Shixiang Gu. A divergence minimization perspective on imitation learning methods. 2019.
- [20] Tianwei Ni, Harshit Sikchi, Yufei Wang, Tejus Gupta, Lisa Lee, and Benjamin Eysenbach. f-irl: Inverse reinforcement learning via state marginal matching. 11 2020.
- [21] Lisa Lee, Benjamin Eysenbach, Emilio Parisotto, Eric P. Xing, Sergey Levine, and Ruslan Salakhutdinov. Efficient exploration via state marginal matching. In *CoRR*, 2019.
- [22] Letian Chen, Rohan R. Paleja, and Matthew C. Gombolay. Learning from suboptimal demonstration via self-supervised reward regression. 2020.
- [23] Ilya Kostrikov, Kumar Krishna Agrawal, Debidatta Dwibedi, Sergey Levine, and Jonathan Tompson. Discriminator-actor-critic: Addressing sample inefficiency and reward bias in adversarial imitation learning, 2018.
- [24] Hana Hoshino, Kei Ota, Asako Kanazaki, and Rio Yokota. Opirl: Sample efficient off-policy inverse reinforcement learning via distribution matching, 2021.

- [25] Wook-Hyun Kwon and S Han. *Receding Horizon Control: Model Predictive Control for State Models*. Springer-Verlag London, 01 2005.
- [26] Yueh-Hua Wu, Nontawat Charoenphakdee, Han Bao, Voot Tangkaratt, and Masashi Sugiyama. Imitation learning from imperfect demonstration. 2019.
- [27] Daniel S. Brown, Wonjoon Goo, and Scott Niekum. Better-than-demonstrator imitation learning via automatically-ranked demonstrations. 2019.
- [28] G. Williams, N. Wagener, B. Goldfain, P. Drews, J. M. Rehg, B. Boots, and E. A. Theodorou. Information theoretic mpc for model-based reinforcement learning. 2017.
- [29] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- [30] David Ha and Jürgen Schmidhuber. World models. *CoRR*, 2018.
- [31] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1861–1870. PMLR, 10–15 Jul 2018.
- [32] Liyiming Ke, Sanjiban Choudhury, Matt Barnes, Wen Sun, Gilwoo Lee, and Siddhartha Srinivasa. Imitation learning as f -divergence minimization, 2021.
- [33] Jonathan Spencer, Sanjiban Choudhury, Arun Venkatraman, Brian Ziebart, and J. Andrew Bagnell. Feedback in imitation learning: The three regimes of covariate shift, 2021.
- [34] Evangelos A. Theodorou and Emanuel Todorov. Relative entropy and free energy dualities: Connections to path integral and kl control. In *2012 IEEE 51st IEEE Conference on Decision and Control (CDC)*, pages 1466–1473, 2012.
- [35] Brian D. Ziebart. *Modeling Purposeful Adaptive Behavior with the Principle of Maximum Causal Entropy*. PhD thesis, Carnegie Mellon University, USA, 2010. AAI3438449.
- [36] Nan Jiang, Alex Kulesza, Satinder P. Singh, and Richard L. Lewis. The dependence of effective planning horizon on model accuracy. In *AAMAS*, 2015.

Checklist

The checklist follows the references. Please read the checklist guidelines carefully for information on how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing the appropriate section of your paper or providing a brief inline description. For example:

- Did you include the license to the code and datasets? **[Yes]** See Section ??.
- Did you include the license to the code and datasets? **[No]** The code and the data are proprietary.
- Did you include the license to the code and datasets? **[N/A]**

Please do not modify the questions and only use the provided macros for your answers. Note that the Checklist section does not count towards the page limit. In your paper, please delete this instructions block and only keep the Checklist section heading above along with the questions/answers below.

1. For all authors...

- (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? **[Yes]**
- (b) Did you describe the limitations of your work? **[Yes]** See Section 5
- (c) Did you discuss any potential negative societal impacts of your work? **[N/A]**
- (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**

2. If you are including theoretical results...

- (a) Did you state the full set of assumptions of all theoretical results? **[Yes]**

- (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#) See Appendix
- 3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[Yes\]](#)
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#)
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[Yes\]](#)
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#)
- 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [\[N/A\]](#)
 - (b) Did you mention the license of the assets? [\[N/A\]](#)
 - (c) Did you include any new assets either in the supplemental material or as a URL? [\[N/A\]](#)
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [\[N/A\]](#)
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[N/A\]](#)
- 5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [\[N/A\]](#)
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [\[N/A\]](#)
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [\[N/A\]](#)

A Gradient Derivation

This section derives the gradient $\frac{\partial \mathcal{L}}{\partial \theta}$ in equation (12) in full details. Recall \mathcal{L} is defined as:

$$\mathcal{L} = \int p_E(V) \log \frac{p_E(V)}{p_Q(V|U_B, \Sigma; \theta)} dV$$

Firstly, we substitute \mathcal{L} in $\frac{\partial \mathcal{L}}{\partial \theta}$ and rewrite $\log \frac{p_E(V)}{p_Q(V|U_B, \Sigma; \theta)}$ as the difference between $\log p_E(V)$ and $p_Q(V|U_B, \Sigma; \theta)$.

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \theta} &= \frac{\partial}{\partial \theta} \int p_E(V) \log \frac{p_E(V)}{p_Q(V|U_B, \Sigma; \theta)} dV \\ &= \int p_E(V) \frac{\partial}{\partial \theta} \log p_E(V) dV - \int p_E(V) \frac{\partial}{\partial \theta} \log p_Q(V|U_B, \Sigma; \theta) dV \\ &= - \int p_E(V) \frac{\partial}{\partial \theta} \log p_Q(V|U_B, \Sigma; \theta) dV \end{aligned} \quad (14)$$

Since $p_E(V)$ is independent of θ , the first derivative in the second line of equation (14) evaluates to 0. Next, we substitute $p_Q(V|U_B, \Sigma; \theta)$ using the optimal control sequence distribution expression in equation (4):

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \theta} &= - \int p_E(V) \frac{\partial}{\partial \theta} \log \frac{1}{Z} p_B(V|U_B, \Sigma) \exp(-\frac{1}{\lambda} S(V; \theta)) dV \\ &= - \int p_E(V) \frac{\partial}{\partial \theta} \log p_B(V|U_B, \Sigma) dV - \int p_E(V) \frac{\partial}{\partial \theta} \log \exp(-\frac{1}{\lambda} S(V; \theta)) dV \\ &\quad + \int p_E(V) \frac{\partial}{\partial \theta} \log Z dV \\ &= - \int p_E(V) \frac{\partial}{\partial \theta} (-\frac{1}{\lambda} S(V; \theta)) dV + \int p_E(V) \frac{\partial}{\partial \theta} \log Z dV \end{aligned} \quad (15)$$

Since $p_B(V|U_B, \Sigma)$ is independent of θ , the first derivative in the second line of equation (14) evaluates to 0. We are left with only two integrals in equation (15).

Next, we factorize out $\frac{\partial}{\partial \theta} \log Z$ from the integral since the partition function Z is constant to all V :

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \theta} &= - \int p_E(V) (-\frac{1}{\lambda} \frac{\partial}{\partial \theta} S(V; \theta)) dV + (\frac{\partial}{\partial \theta} \log Z) \int p_E(V) dV \\ &= - \int p_E(V) (-\frac{1}{\lambda} \frac{\partial}{\partial \theta} S(V; \theta)) dV + \frac{1}{Z} \frac{\partial Z}{\partial \theta} \end{aligned} \quad (16)$$

The second line in equation (16) follows as $\int p_E(V) dV = 1$.

Next, we substitute $Z = \int p_B(V|U_B, \Sigma) \exp(-\frac{1}{\lambda} S(V; \theta)) dV$ in $\frac{\partial Z}{\partial \theta}$ and simplify it:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \theta} &= - \int p_E(V) (-\frac{1}{\lambda} \frac{\partial}{\partial \theta} S(V; \theta)) dV + \frac{1}{Z} \frac{\partial}{\partial \theta} \int p_B(V|U_B, \Sigma) \exp(-\frac{1}{\lambda} S(V; \theta)) dV \\ &= - \int p_E(V) (-\frac{1}{\lambda} \frac{\partial}{\partial \theta} S(V; \theta)) dV + \frac{1}{Z} \int p_B(V|U_B, \Sigma) \frac{\partial}{\partial \theta} \exp(-\frac{1}{\lambda} S(V; \theta)) dV \\ &= - \int p_E(V) (-\frac{1}{\lambda} \frac{\partial}{\partial \theta} S(V; \theta)) dV \\ &\quad + \frac{1}{Z} \int p_B(V|U_B, \Sigma) \exp(-\frac{1}{\lambda} S(V; \theta)) \frac{\partial}{\partial \theta} (-\frac{1}{\lambda} S(V; \theta)) dV \\ &= - \int p_E(V) (-\frac{1}{\lambda} \frac{\partial}{\partial \theta} S(V; \theta)) dV \\ &\quad + \int \frac{p_B(V|U_B, \Sigma) \exp(-\frac{1}{\lambda} S(V; \theta))}{Z} \frac{\partial}{\partial \theta} (-\frac{1}{\lambda} S(V; \theta)) dV \end{aligned} \quad (17)$$

We rewrite $\frac{1}{\lambda} p_B(V|U_B, \Sigma) \exp(-\frac{1}{\lambda} S(V; \theta))$ in the last line of equation (17) as $p_Q(V|U_B, \Sigma; \theta)$ using equation (4) and finally we have:

$$\frac{\partial \mathcal{L}}{\partial \theta} = \int p_E(V) \left(\frac{1}{\lambda} \frac{\partial}{\partial \theta} S(V; \theta) \right) dV - \int p_Q(V|U_B, \Sigma; \theta) \left(\frac{1}{\lambda} \frac{\partial}{\partial \theta} S(V; \theta) \right) dV \quad (18)$$

B Proof of Theorem 3.1

We use T to denote the task horizon, and K to denote the receding time horizon for local optimization. Let us simplify notations in the optimal control sequence distribution $p_Q(V|U_B, \Sigma, x_t; \theta)$ and remove the explicit dependency on U_B and Σ . We assume that all control sequences are applied with U_B as the base distribution and Σ as the covariate matrix. We assume the expert's underlying cost function is parameterized by θ_E , so we have $p_E(\cdot) = p(\cdot; \theta_E)$.

B.1 Sketch

We present a theoretical analysis on the convergence of RHIRL. Our main theorem 3.1 states that given that the Kullback–Leibler (KL) divergence over local control sequence distribution for each time step $t = 0, 1, \dots, T-1$ is bounded by ϵ , though we do not query expert during the learning, using the cached expert demonstrations alone allows us to bound the error over global state marginal distribution *linear* in the task horizon T under total variance measure.

First, we show that if the KL-divergence over the local control sequence distribution is bounded by ϵ , so is the KL-divergence over the resulting state distribution. At each time step t , the optimal control sequences distribution $p_Q(V_t|x_t; \theta_t)$ at the initial state x_t contains the full information to generate the corresponding state trajectories $p(\tau_t|x_t; \theta)$, and consequently the state distributions $p_t(x|x_t; \theta)$ (by neglecting the temporal information). Upon applying the information loss (lemma B.1), we prove in lemma B.2 that, given initial state x_t , if the KL-divergence over V_t is bound by ϵ , so is the corresponding state distribution $p_t(x)$, i.e. $D_{\text{KL}}(p_t(x|x_t; \theta_E) \parallel p_t(x|x_t; \theta)) \leq D_{\text{KL}}(p(V_t|x_t; \theta_E) \parallel p_Q(V_t|x_t; \theta)) \leq \epsilon$.

Next, we show the KL-divergence over global state marginal distribution between two consecutive time steps is also bounded by ϵ . For each optimal control sequence we compute, we only execute the first control and use the rest to warm start the re-planning for the next time step. Therefore, for each time step, we only change a small region of the global state distribution, i.e. reachable space of the current time step. We use $p_{\text{RHC}}^t(x; \theta)$ to denote the global state marginal distribution by recursively applying RHC from time step $i = 0, \dots, t$ under θ and switching to θ_E thereafter until $T-1$. Using generalized log sum inequality, we prove in lemma B.3 that if the KL-divergence over V_t is bounded by ϵ for all $t = 0, \dots, T-1$, the KL-divergence over the global state marginal distribution between each of the two consecutive time steps is bounded by ϵ , i.e. $D_{\text{KL}}(p_{\text{RHC}}^t(x; \theta) \parallel p_{\text{RHC}}^{t+1}(x; \theta)) \leq \frac{K+1}{T} \epsilon$.

Finally, we use Pinsker's inequality to upper bound the total variation (TV) distance by KL-divergence over state marginal distribution. Then we use the triangle inequality to show that the TV distance between expert and the actual visited state distribution over the task horizon T using RHIRL is bounded by an error linearly in T , i.e. $D_{\text{TV}}(p_E(x) \parallel p_{\text{RHC}}(x; \theta)) < T\sqrt{\epsilon/2}$.

B.2 Proofs

First, we use **Lemma B.1, B.2** to prove that the KL-divergence over control sequence space upper bounds the KL-divergence over the resulting state distribution. We define the control sequence starts at task time step t as $V_t = \{v_t, v_{t+1}, \dots, v_{t+K-1}\}$. Moreover, its corresponding trajectory segment $\tau_t = \{x_t, x_{t+1}, \dots, x_{t+K}\}$ is computed uniquely from V_t and initial state x_t by iteratively applying the dynamic model $x_{t+1} = f(x_t, v_t)$. We use $p(x_t)$ to denote the state density at a single time step t . Assume V_t is optimized based on the cost parameterization θ , then the corresponding state distribution is defined as the summation of all state density over horizon K , i.e. $p_t(x; \theta) = \frac{1}{K+1} \sum_{i=t}^{t+K} p(x_i; \theta)$.

Lemma B.1 (Information loss[32]). *Let a and b be two random variables and $f(\cdot)$ be a convex function. Let $P(a, b)$ be a joint probability distribution. The marginal distributions are $P(a) = \sum_b P(a, b)$ and $P(b) = \sum_a P(a, b)$. Assume that a can explain away b . This is expressed as follows – given any two probability distribution $P(\cdot)$, $Q(\cdot)$, assume the following equality holds for all a, b :*

$$P(b|a) = Q(b|a) \quad (19)$$

Under these conditions, the following inequality holds:

$$\sum_a Q(a) f\left(\frac{P(a)}{Q(a)}\right) > \sum_b Q(b) f\left(\frac{P(b)}{Q(b)}\right) \quad (20)$$

Lemma B.2. Given the initial state x_t and the two control sequence distributions $p(V_t|x_t; \theta_E)$ and $p_Q(V_t|x_t; \theta)$, the KL-divergence between the resulting state distribution is upper bounded by KL-divergence between the control sequence distribution.

$$D_{\text{KL}}(p_t(x|x_t; \theta_E) \parallel p_t(x|x_t; \theta)) \leq D_{\text{KL}}(p(V_t|x_t; \theta_E) \parallel p_Q(V_t|x_t; \theta)) \quad (21)$$

Proof. Firstly, we prove that the KL-divergence over state trajectory distribution $p(\tau_t)$ is upper bounded by the KL-divergence between $p(V_t)$. Given dynamical model f , the control sequence V_t and the initial state x_t contains all information to generate the corresponding τ_t . Therefore, for any joint distribution $P(\tau_t, V_t|x_t)$ and $Q(\tau_t, V_t|x_t)$, the following is true

$$P(\tau_t|V_t, x_t) = Q(\tau_t|V_t, x_t)$$

Upon applying the information loss Lemma B.1, we have the inequality:

$$D_{\text{KL}}(p(\tau_t|x_t; \theta_E) \parallel p(\tau_t|x_t; \theta)) \leq D_{\text{KL}}(p(V_t|x_t; \theta_E) \parallel p_Q(V_t|x_t; \theta)) \quad (22)$$

Next we prove that the KL-divergence between state distribution is upper bounded by the trajectory distribution. Since a trajectory $\tau_t = \{x_t, x_{t+1}, \dots, x_{t+K}\}$ contains full information of the resulting states (by neglecting the temporal information), for any joint distribution $P(x|\tau_t)$ and $Q(x|\tau_t)$, the following is true

$$P(x|\tau_t) = Q(x|\tau_t) \quad (23)$$

Upon applying Lemma B.1 we have the inequality:

$$D_{\text{KL}}(p_t(x|x_t; \theta_E) \parallel p_t(x|x_t; \theta)) \leq D_{\text{KL}}(p(\tau_t|x_t; \theta_E) \parallel p(\tau_t|x_t; \theta)) \quad (24)$$

Therefore, given the KL-divergence between the control sequence distribution is upper bounded by ϵ , we use the equality in equation (22) and (24) to show that the KL-divergence between the resulting state distribution is also upper bounded by ϵ .

$$D_{\text{KL}}(p_t(x|x_t; \theta_E) \parallel p_t(x|x_t; \theta)) \leq D_{\text{KL}}(p(V_t|x_t; \theta_E) \parallel p_Q(V_t|x_t; \theta)) \leq \epsilon \quad (25)$$

□

Definition 1 (One-step Recoverability [33]). Assume that the state distribution of the learner and expert are different at time t , that is $D_{\text{KL}}(p(x_t; \theta_E) \parallel p(x_t; \theta)) \neq 0$, there exists a policy π_{re} that when used for the learner, can bound:

$$D_{\text{KL}}(p(x_{t+1}; \theta_E) \parallel p(x_{t+1}; \pi_{re})) \leq \epsilon_1 \quad (26)$$

where the current initial state distribution of the student follows $p(x_t; \theta)$.

Intuitively, this condition requires that, no matter what is the current state distribution, the learner can recover to the expert demonstrated distribution in a single time-step. In our case, this is a natural condition since the difference in the initial state distribution $p(x_t; \theta)$ and $p(x_t; \theta_E)$ is not arbitrarily large: we use re-planning to ensure the receding state sequences is always bounded below ϵ , hence this recoverability condition can be easily satisfied. We emphasize that this recoverable policy is **never** executed in our algorithm, it is only used for the theoretical analysis. Next, we derive a bound over the global state marginal distribution between two consecutive time steps. At each time step $t = 0, \dots, T-1$, we re-optimize the local control sequence distribution and only execute the first control, hence we only change state density over a small reachable space. We define $p_{\text{RHC}}^t(x, \theta)$ as the global marginal state distribution by applying RHC from $i = 0, 1, \dots, t$ under θ and then using the recoverable policy π_{re} to switch to the expert θ_E thereafter until $T-1$, i.e., $p_{\text{RHC}}^t(x; \theta) = \frac{1}{T}(\sum_{i=0}^{t+1} p(x_i; \theta) + p(x_{t+2}; \pi_{re}) + \sum_{i=t+3}^T p(x_i; \theta_E))$. According to the definition of the recoverable policy, we have $D_{\text{KL}}(p(x_{t+2}; \theta_E) \parallel p(x_{t+2}; \pi_{re})) \leq \epsilon_1$, therefore, $p_{\text{RHC}}^t(x; \theta) \approx \frac{1}{T}(\sum_{i=0}^{t+1} p(x_i; \theta) + \sum_{i=t+2}^T p(x_i; \theta_E))$. To quantify the change in global state marginal distribution, we derive a bound for the KL-divergence between two consecutive time steps, i.e. $D_{\text{KL}}(p_{\text{RHC}}^{t-1}(x; \theta) \parallel p_{\text{RHC}}^t(x; \theta))$.

Lemma B.3. *If the KL-divergence over resulting state density from the control sequence distribution of length K are bounded by ϵ , i.e. $D_{\text{KL}}(p_t(x|x_t; \theta_E) \parallel p_t(x|x_t; \theta)) < \epsilon$, where x_t is the state encountered by our policy at $t = 0, 1, \dots, T-1$ and is one-step recoverable, then KL-divergence over the global state marginal distribution between two consecutive control executions are bounded by ϵ ,*

$$D_{\text{KL}}(p_{\text{RHC}}^{t-1}(x; \theta) \parallel p_{\text{RHC}}^t(x; \theta)) \leq \frac{K+1}{T} \epsilon \quad (27)$$

for $t = 1, \dots, T-1$.

We state the generalized log sum inequality below in lemma B.4, the proof can be found in the Appendix of [32]. Lemma B.4 and B.5 will be used in the proof for Lemma B.3.

Lemma B.4 (Generalized log sum inequality[32]). *Let p_1, \dots, p_n and q_1, \dots, q_n be non-negative numbers. Let $p = \sum_{i=1}^n p_i$ and $q = \sum_{i=1}^n q_i$. Let $f(\cdot)$ be a convex function. We have the following:*

$$\sum_{i=1}^n q_i f\left(\frac{p_i}{q_i}\right) \geq q f\left(\frac{p}{q}\right) \quad (28)$$

Lemma B.5. *Let $p(x)$ and $q(x)$ be non-negative functions, and c is a constant factor. We have the following:*

$$\int c p(x) \log \frac{c p(x)}{c q(x)} = c \int p(x) \log \frac{p(x)}{q(x)} \quad (29)$$

Proof.

$$\int c p(x) \log \frac{c p(x)}{c q(x)} dx = \int c p(x) \log \frac{p(x)}{q(x)} dx = c \int p(x) \log \frac{p(x)}{q(x)} dx \quad (30)$$

□

Now, we are ready to prove lemma B.3.

Proof. For each $t = 0, \dots, T-1$, we re-plan for the optimal local control sequence start at x_t so that the resulting state distribution over horizon K is bounded, i.e. $D_{\text{KL}}(p_t(x|x_t; \theta_E) \parallel p_t(x|x_t; \theta)) < \epsilon$, where $p_t(x|x_t; \theta) = \frac{1}{K+1} \sum_{i=t}^{t+K} p(x_i; \theta)$. However, instead of executing all K controls in the sequence, we only execute the first control at the current time step t and change the state distribution $p(x_{t+1}; \theta)$ reachable for that single time step, then we use the remaining control sequence to warm start the local control sequence optimization for the next time step. To account the effect of replanning, for each time step t , since we do not change the state distribution after $p(x_{t+1}; \theta)$, we can think of the change in the global state distribution as if we follow the optimal control under θ at time step t and then use the recoverable policy π_{re} to switch to θ_E afterwards over the control sequence horizon K . Hence, the actual state density is $\frac{1}{K+1} (p(x_t; \theta) + p(x_{t+1}; \theta) + p(x_{t+2}; \pi_{re}) + \sum_{i=t+3}^{t+K} p(x_i; \theta_E))$. Theoretically, since we do not query the expert online, the initial state x_t distribution in theorem 3.1 and lemma B.2 should follow the expert demonstration at time t , i.e. $p(x_t; \theta_E)$. However, our MPC controller cannot jump to this distribution and we replan from our current state distribution $p(x_t; \theta)$. To resolve this mismatch, we require the recoverability condition in our optimization procedure such that, for each resulting state from the controller, there always exists a one-step recoverable policy π_{re} that can correct the current state distribution $p(x_t; \theta)$ to $p(x_{t+1}; \theta_E)$ in one step. Therefore, with this one-step recoverable condition on every state x_t induced by the cost function parameterized by θ , x_t in theorem 3.1 and lemma B.2 now follows the state distribution of our controller, i.e. $p(x_t; \theta)$.

That is,

$$\begin{aligned} & D_{\text{KL}}(p_t(x|x_t; \theta_E) \parallel \frac{1}{K+1} (p(x_t) + p(x_{t+1}; \theta) + \sum_{i=t+2}^{t+K} p(x_i; \theta_E))) \\ &= D_{\text{KL}}(\frac{1}{K+1} (p(x_t; \theta) + \sum_{i=t+1}^{t+K} p(x_i; \theta_E)) \parallel \frac{1}{K+1} (p(x_t; \theta) + p(x_{t+1}; \theta) + \sum_{i=t+2}^{t+K} p(x_i; \theta_E))) \\ &\leq D_{\text{KL}}(\frac{1}{K+1} (p(x_t; \theta) + \sum_{i=t+1}^{t+K} p(x_i; \theta_E)) \parallel \frac{1}{K+1} (p(x_t; \theta) + \sum_{i=t+1}^{t+K} p(x_i; \theta))) \\ &= D_{\text{KL}}(p_t(x|x_t; \theta_E) \parallel p_t(x|x_t; \theta)) \leq \epsilon \end{aligned} \quad (31)$$

Recall that we use $p_{\text{RHC}}^t(x; \theta)$ to account for the global state marginal distribution resulted from executing a single optimal control at time step t . More specifically, $p_{\text{RHC}}^t(x; \theta)$ is defined as the state marginal distribution by executing only the first optimal control from the replanned optimal control sequence at each time step from $i = 0, \dots, t$ under θ and switching to θ_E thereafter by using the recoverable policy π_{re} until $T - 1$, i.e. $p_{\text{RHC}}^t(x; \theta) \approx \frac{1}{T}(\sum_{i=0}^{t+1} p(x_i; \theta) + \sum_{i=t+2}^T p(x_i; \theta_E))$. We bound the KL-divergence over global state distribution between two consecutive time step as follow:

$$\begin{aligned}
& D_{\text{KL}}(p_{\text{RHC}}^{t-1}(x; \theta) \parallel p_{\text{RHC}}^t(x; \theta)) \\
&= D_{\text{KL}}\left(\frac{1}{T}\left(\sum_{i=0}^t p(x_i; \theta) + p(x_{t+1}; \pi_{re}) + \sum_{i=t+2}^T p(x_i; \theta_E)\right) \parallel \frac{1}{T}\left(\sum_{i=0}^{t+1} p(x_i; \theta) + p(x_{t+2}; \pi_{re}) + \sum_{i=t+3}^T p(x_i; \theta_E)\right)\right) \\
&\approx D_{\text{KL}}\left(\frac{1}{T}\left(\sum_{i=0}^t p(x_i; \theta) + \sum_{i=t+1}^T p(x_i; \theta_E)\right) \parallel \frac{1}{T}\left(\sum_{i=0}^{t+1} p(x_i; \theta) + \sum_{i=t+2}^T p(x_i; \theta_E)\right)\right) \\
&= \frac{1}{T} \int \left(\sum_{i=0}^t p(x_i; \theta) + \sum_{i=t+1}^T p(x_i; \theta_E)\right) \log \frac{\sum_{i=0}^t p(x_i; \theta) + \sum_{i=t+1}^T p(x_i; \theta_E)}{\sum_{i=0}^{t+1} p(x_i; \theta) + \sum_{i=t+2}^T p(x_i; \theta_E)} dx \\
&\leq \frac{1}{T} \left(\int \sum_{i=0}^{t-1} p(x_i; \theta) \log \frac{\sum_{i=0}^{t-1} p(x_i; \theta)}{\sum_{i=0}^{t-1} p(x_i; \theta)} dx \right. \\
&\quad + \int \left(p(x_t; \theta) + \sum_{i=t+1}^{t+K} p(x_i; \theta_E)\right) \log \frac{p(x_t; \theta) + \sum_{i=t+1}^{t+K} p(x_i; \theta_E)}{p(x_t; \theta) + p(x_{t+1}; \theta) + \sum_{i=t+2}^{t+K} p(x_i; \theta_E)} dx \\
&\quad \left. + \int \sum_{i=t+K+1}^T p(x_i; \theta_E) \log \frac{\sum_{i=t+K+1}^T p(x_i; \theta_E)}{\sum_{i=t+K+1}^T p(x_i; \theta_E)} dx \right) \\
&= \frac{1}{T} \int \left(p(x_t; \theta) + \sum_{i=t+1}^{t+K} p(x_i; \theta_E)\right) \log \frac{p(x_t; \theta) + \sum_{i=t+1}^{t+K} p(x_i; \theta_E)}{p(x_t; \theta) + p(x_{t+1}; \theta) + \sum_{i=t+2}^{t+K} p(x_i; \theta_E)} dx \\
&= \frac{K+1}{T} \int \frac{1}{K+1} \left(p(x_t; \theta) + \sum_{i=t+1}^{t+K} p(x_i; \theta_E)\right) \log \frac{\frac{1}{K+1} \left(p(x_t; \theta) + \sum_{i=t+1}^{t+K} p(x_i; \theta_E)\right)}{\frac{1}{K+1} \left(p(x_t; \theta) + p(x_{t+1}; \theta) + \sum_{i=t+2}^{t+K} p(x_i; \theta_E)\right)} dx \\
&= \frac{K+1}{T} D_{\text{KL}}\left(\frac{1}{K+1} \left(p(x_t; \theta) + \sum_{i=t+1}^{t+K} p(x_i; \theta_E)\right) \parallel \frac{1}{K+1} \left(p(x_t; \theta) + p(x_{t+1}; \theta) + \sum_{i=t+2}^{t+K} p(x_i; \theta_E)\right)\right) \\
&= \frac{K+1}{T} D_{\text{KL}}(p_t(x|x_t; \theta_E) \parallel \frac{1}{K+1} \left(p(x_t; \theta) + p(x_{t+1}; \theta) + \sum_{i=t+2}^{t+K} p(x_i; \theta_E)\right)) \\
&\leq \frac{K+1}{T} \epsilon
\end{aligned} \tag{32}$$

The first equality in equation (32) follows the definition of $p_{\text{RHC}}^t(x; \theta)$ and the second line follows the definition of the recoverable policy π_{re} . Then, we use lemma B.5 to factor out $\frac{1}{T}$ in the third line. The next inequality follows from the generalized log sum inequality stated in lemma B.4, and we have the first and third terms reduce to 0 and are left with the second term in the next line. We apply lemma B.5 again to the integral using the constant factor $\frac{1}{K+1}$. In addition, to make the equality hold, we multiply the inverse of the constant factor $K+1$ outside the integral. We observe the integral is now the KL divergence between the expert $p_t(x|x_t; \theta_E)$ and one-step-execution of our policy $\frac{1}{K+1} (p(x_t; \theta) + p(x_{t+1}; \theta) + \sum_{i=t+2}^{t+K} p(x_i; \theta_E))$. The final inequality follows from the bound derived in equation (31).

For $t = 0$, we have $p_{\text{RHC}}^0(x; \theta) = \frac{1}{T}(p(x_0; \theta) + \sum_{i=1}^{T-1} p(x_i; \theta_E))$. Since the initial state x_0 for expert and our policy are sampled from the same initial state distribution μ , $p(x_0)$ is independent of θ , i.e. $p(x_0; \theta) = p(x_0; \theta_E)$. Therefore, $p_{\text{RHC}}(x; \theta_E) = p(x_0) + \sum_{i=1}^{T-1} p(x_i; \theta_E) = p_{\text{RHC}}^0(x; \theta)$. Moreover, the final global state marginal distribution $p_{\text{RHC}}(x; \theta)$ is the same as the $p_{\text{RHC}}^{T-1}(x; \theta)$,

i.e. $p_{\text{RHC}}(x; \theta) = \sum_{i=0}^{T-1} p(x_i; \theta) = p_{\text{RHC}}^{T-1}(x; \theta)$. For any $t = 1, 2, \dots, T-1$, we have proved $D_{\text{KL}}(p_{\text{RHC}}^{t-1}(x; \theta) \parallel p_{\text{RHC}}^t(x; \theta)) \leq \frac{K+1}{T}\epsilon$.

□

Finally, we are prepared to prove theorem 3.1.

Proof. We evaluate the TV distance over the state marginal distribution between the expert policy and our control law.

$$\begin{aligned} D_{\text{TV}}(p_{\text{RHC}}(x; \theta_E) \parallel p_{\text{RHC}}(x; \theta)) &= D_{\text{TV}}(p_{\text{RHC}}^0(x; \theta) \parallel p_{\text{RHC}}^{T-1}(x; \theta)) \\ &\leq \sum_{t=1}^{T-1} D_{\text{TV}}(p_{\text{RHC}}^{t-1}(x; \theta) \parallel p_{\text{RHC}}^t(x; \theta)) \end{aligned} \quad (33)$$

The first equality in equation (33) follows from the fact that $p_{\text{RHC}}(x; \theta_E) = p_{\text{RHC}}^0(x; \theta)$ and $p_{\text{RHC}}(x; \theta) = p_{\text{RHC}}^{T-1}(x; \theta)$. We use triangle inequality of the TV distance measures to obtain the inequality in the second line.

Recall that by Pinsker’s inequality, the total variation (TV) distance is related to Kullback–Leibler (KL) divergence by the following inequality: $D_{\text{TV}}(P \parallel Q) \leq \sqrt{\frac{1}{2} D_{\text{KL}}(P \parallel Q)}$. We apply Pinsker’s inequality to each of the TV terms in the second line of equation (33) to bound them by a summation of KL-divergence as shown in the first line of equation (34). Next, given the control sequence distribution for every time step is bounded by ϵ , we apply lemma B.2 to show that the resulting state distribution from the optimal control sequences for each time step t is also bounded by ϵ . Next, we use this result and apply Lemma B.3 to bound the KL-divergence over the global state marginal distribution between two consecutive time steps by $\frac{K+1}{T}\epsilon$. Second line in equation (34) follows from this result and finally we derive the final bound linear in T .

$$\begin{aligned} D_{\text{TV}}(p_{\text{RHC}}(x; \theta_E) \parallel p_{\text{RHC}}(x; \theta)) &\leq \sum_{t=1}^{T-1} \sqrt{\frac{1}{2} D_{\text{KL}}(p_{\text{RHC}}^{t-1}(x; \theta) \parallel p_{\text{RHC}}^t(x; \theta))} \\ &\leq \sum_{t=1}^{T-1} \sqrt{\frac{(K+1)\epsilon}{2T}} \\ &= (T-1) \sqrt{\frac{K+1}{T}} \sqrt{\epsilon/2} \\ &\leq T \sqrt{\epsilon/2} \end{aligned} \quad (34)$$

The last line follows from the fact that $K \ll T$, so $\sqrt{\frac{K+1}{T}} < 1$.

□

C Extension to Stochastic Dynamics

RHRL optimizes the trajectories in the space of control sequences $p(V)$, whereas $V = \{u_0, u_1, u_2, \dots, u_{K-1}\}$ is a sequence of controls. If the system is deterministic, we can apply V to the dynamical system $x_{t+1} = f(x_t, v_t)$ with start state x_0 and obtain a state sequence $\tau = (x_0, x_1, \dots, x_K)$. We recall that total state trajectory cost of V defined in equation (3) as follows:

$$S(V, x_0; \theta) = \sum_{k=0}^K g(x_k; \theta)$$

We use the information-theoretic MPC (MPPI) [28] to solve for an optimal control sequence distribution at the current start state x_t in iteration t . The main result of MPPI suggests that, under a deterministic system, the optimal control sequence distribution Q minimizes the “free energy” of the dynamical system and this free energy can be calculated from the cost of the state trajectory under Q . Mathematically, the probability density $p_Q(V^*)$, as shown in equation (4) can be expressed as a

function of the state cost $S(V, x_t; \theta)$, with respect to a Gaussian “base” distribution $B(V_B, \Sigma)$ that depends on the control cost:

$$p_Q(V^*|U_B, \Sigma, x_t; \theta) = \frac{1}{Z} p_B(V^*|U_B, \Sigma) \exp\left(-\frac{1}{\lambda} S(V^*, x_t; \theta)\right),$$

where Z is the partition function. Intuitively, this result shows that the control sequence V that results in lower state-trajectory cost $S(V)$ are exponentially more likely to be chosen.

In this section, we extend RHIRL to stochastic dynamics where $x_{t+1} \sim p(x_{t+1}|x_t, v_t)$. Due to the stochasticity of the dynamics, the state trajectory cost in equation (3) and the optimal control distribution in equation (4) are affected. Therefore, we first redefine the trajectory state cost under stochastic dynamics, then derive the counterpart of equation (4) for the optimal control sequence distribution under the stochastic dynamics, finally we adapt our existing RHIRL algorithm to stochastic dynamics.

C.1 State Trajectory Cost $\tilde{S}(V; \theta)$ under Stochastic Dynamics

Due to the stochasticity of the dynamics, given the initial state x_0 , we no longer have a one-to-one mapping from the control sequence V to the resulting state trajectory $\tau = (x_0, x_1, \dots, x_K)$. Instead, we have a distribution of state trajectories:

$$p(\tau|x_0, V) = \prod_{t=0}^{K-1} p(x_{t+1}|x_t, v_t) \quad (35)$$

To accommodate this change, the trajectory state cost of a control sequence $\tilde{S}(V, x_0; \theta)$ is defined over the distribution of the resulting state trajectories, instead of single trajectory:

$$\tilde{S}(V, x_0; \theta) = \int p(\tau|x_0, V) S(\tau|x_0; \theta) d\tau \quad (36)$$

$$= \int \prod_{t=0}^{K-1} p(x_{t+1}|x_t, v_t) \sum_{t=0}^K g(x_t; \theta) d\tau \quad (37)$$

We always measure the preferences over the control sequence V by their resulting state trajectories τ . Hence, when the resulting trajectories changes from a single deterministic sequence of states to a distribution of state trajectories, we adapt our measure of the resulting cost: under the deterministic dynamics where each V uniquely maps to the same state trajectory τ , the state trajectory cost is the cost of that specific trajectory; while under the stochastic dynamics where the same control sequence V maps to a distribution of τ , the state trajectory cost of a control sequence is now defined as the expected state cost of the distribution of trajectory. We measure the state trajectory cost of a control sequence $\tilde{S}(V, x_0; \theta)$, instead of simply a state trajectory cost on the states itself $S(\tau, x_0; \theta)$, because we want to use this measure to directly optimize the control sequence.

C.2 Optimal Control Sequence Distribution under Stochastic Dynamics

Next, we derive the optimal control sequence distribution under the stochastic dynamics. Our derivation is based on MPPI [28], which uses the “free-energy” principle to derive the optimal control sequence distribution under deterministic dynamics.

Definition 2 (Free-energy ([34], Definition 1)). Let $\mathbb{P} \in \mathcal{P}(\mathcal{Z})$ and the function $\mathcal{J}(x): \mathcal{Z} \rightarrow \mathbb{R}$ be a measurable function. The the term:

$$\mathbb{E}(\mathcal{J}(x)) = \log \int \exp(\rho \mathcal{J}(x)) d\mathbb{P} \quad (38)$$

is called free energy of $\mathcal{J}(x)$ with respect to \mathbb{P} , ρ is a constant.

Now we have the free-energy of a control system under stochastic dynamics as stated below. It has a “Gaussian” base control sequence distribution $B(U_B, \Sigma)$ such that its control sequence distribution follows $p_B(V^*|U_B, \Sigma)$ whereas Σ is the Gaussian control noise covariance matrix. $\tilde{S}(V; \theta)$ denotes the state trajectory cost function.

$$\mathcal{F}(\tilde{S}, p_B, x_0, \lambda; \theta) = \log(\mathbb{E}_{p_B}[\exp(-\frac{1}{\lambda} \tilde{S}(V, x_0; \theta))]), \quad (39)$$

$\lambda \in \mathbb{R}^+$ is the inverse temperature of the control system.

Suppose now we have another control sequence distribution with probability measure $p(V)$ and these two distributions are absolutely continuous, then we can rewrite the free-energy w.r.t $p_B(V)$ using the expectation over the density of $p(V)$ use the standard importance sampling trick:

$$\mathcal{F}(\tilde{S}, p_B, x_0, \lambda; \theta) = \log(\mathbb{E}_{p_B}[\exp(-\frac{1}{\lambda}\tilde{S}(V, x_0; \theta))]) \quad (40)$$

$$= \log(\mathbb{E}_{p_B}[\exp(-\frac{1}{\lambda}\tilde{S}(V, x_0; \theta) \frac{p_B(V^*|U_B, \Sigma)}{p(V)})]) \quad (41)$$

$$\geq \mathbb{E}_p[\log(\exp(-\frac{1}{\lambda}\tilde{S}(V, x_0; \theta) \frac{p_B(V^*|U_B, \Sigma)}{p(V)}))] \quad (42)$$

$$= -\frac{1}{\lambda} \mathbb{E}_p[\tilde{S}(V, x_0; \theta) + \lambda \log(\frac{p(V)}{p_B(V^*|U_B, \Sigma)})] \quad (43)$$

$$= -\frac{1}{\lambda} (\mathbb{E}_p[\tilde{S}(V, x_0; \theta)] + \lambda \mathbb{E}_p[\log(\frac{p(V)}{p_B(V^*|U_B, \Sigma)})]) \quad (44)$$

$$= -\frac{1}{\lambda} (\mathbb{E}_p[\tilde{S}(V, x_0; \theta)] + \lambda D_{\text{KL}}(p(V) || p_B(V^*|U_B, \Sigma))) \quad (45)$$

Therefore, we have

$$\mathcal{F}(S, p_B, x_0, \lambda; \theta) \geq -\frac{1}{\lambda} (\mathbb{E}_p[\tilde{S}(V, x_0; \theta)] + \lambda D_{\text{KL}}(p(V) || p_B(V^*|U_B, \Sigma))) \quad (46)$$

The right-hand side is the lower bound of the free-energy of the control system. We use $p_Q(V)$ to denote the optimal control sequence distribution. This distribution is only optimal if and only if the bound in the equation above is tight, i.e. $\mathcal{F}(\tilde{S}, p_B, x_0, \lambda; \theta) = -\frac{1}{\lambda} \mathbb{E}_{p_Q}[\tilde{S}(V, x_0; \theta)] + \lambda D_{\text{KL}}(p_Q(V) || p_B(V^*|U_B, \Sigma))$.

We claim that the optimal control sequence distribution \tilde{p}_Q under the stochastic dynamics is as follows:

$$\tilde{p}_Q(V^*|U_B, \Sigma, x_0; \theta) = \frac{1}{Z} \exp(-\frac{1}{\lambda}\tilde{S}(V, x_0; \theta)) p_B(V^*|U_B, \Sigma), \quad (47)$$

whereas $Z = \int \exp(-\frac{1}{\lambda}\tilde{S}(V, x_0; \theta)) p_B(V^*|U_B, \Sigma) dV$ is the partition function.

We prove that equation (47) is the optimal control sequence distribution under stochastic dynamics by showing that this $\tilde{p}_Q(V^*|U_B, \Sigma, x_0; \theta)$ tightens the bound of free-energy in equation (46). We substitute $\tilde{p}_Q(V^*|U_B, \Sigma, x_0; \theta)$ into the RHS of equation (46) and simplify the expression of the KL-divergence:

$$\mathcal{F}(\tilde{S}, p_B, x_0, \lambda; \theta) \geq -\frac{1}{\lambda} (\mathbb{E}_{\tilde{p}_Q}[\tilde{S}(V, x_0; \theta)] + \lambda D_{\text{KL}}(\tilde{p}_Q(V^*|U_B, \Sigma, x_0; \theta) || p_B(V^*|U_B, \Sigma))) \quad (48)$$

$$= -\frac{1}{\lambda} (\mathbb{E}_{\tilde{p}_Q}[\tilde{S}(V, x_0; \theta)] + \lambda \mathbb{E}_{\tilde{p}_Q}[\log(\frac{\tilde{p}_Q(V^*|U_B, \Sigma, x_0; \theta)}{p_B(V^*|U_B, \Sigma)})]) \quad (49)$$

$$= -\frac{1}{\lambda} \mathbb{E}_{\tilde{p}_Q}[\tilde{S}(V, x_0; \theta)] - \mathbb{E}_{\tilde{p}_Q}[\log(\frac{\frac{1}{Z} \exp(-\frac{1}{\lambda}\tilde{S}(V, x_0; \theta)) p_B(V^*|U_B, \Sigma)}{p_B(V^*|U_B, \Sigma)})] \quad (50)$$

$$= -\frac{1}{\lambda} \mathbb{E}_{\tilde{p}_Q}[\tilde{S}(V, x_0; \theta)] - (\frac{1}{\lambda} \mathbb{E}_{\tilde{p}_Q}[\tilde{S}(V, x_0; \theta)] - \log(Z)) \quad (51)$$

Next we substitute the expression for the partition function Z and we found that the RHS is exactly the definition of the the free-energy of the control system with base distribution $B(U_B, \Sigma)$:

$$\mathcal{F}(\tilde{S}, p_B, x_0, \lambda; \theta) \geq \log(Z) \quad (52)$$

$$= \log\left(\int \exp\left(-\frac{1}{\lambda}\tilde{S}(V, x_0; \theta)\right)p_B(V^*|U_B, \Sigma)dV\right) \quad (53)$$

$$= \log(\mathbb{E}_{p_B}[\exp(-\frac{1}{\lambda}\tilde{S}(V, x_0; \theta))]) \quad (54)$$

$$= \mathcal{F}(\tilde{S}, p_B, x_0, \lambda; \theta) \quad (55)$$

The final equality forces the inequality to be tight. Therefore, $p_Q(V^*|U_B, \Sigma, x_0; \theta)$ in equation (47) is the optimal control sequence distribution under the stochastic dynamics.

We observe that the optimal control sequence distribution under deterministic system $p_Q(V^*|U_B, \Sigma, x_0; \theta)$ in equation (4) and that under the stochastic dynamics $\tilde{p}_Q(V^*|U_B, \Sigma, x_0; \theta)$ in equation (47) only differs in the calculation of the state trajectory cost of the control sequences. Intuitively, it means that under deterministic dynamics, we choose the control sequence V that will, for sure, leads to a state trajectory with lower cost; while when we extend to stochastic dynamics, the control sequence V that results in lower state-trajectory cost $S(V)$ in expectation are exponentially more likely to be chosen. Practically, now we need more samples for a single control sequence to compute the expectation in equation (37).

C.3 RHIRL under Stochastic Dynamics

Next, we adapt our RHIRL algorithm to this new state trajectory cost measure $\tilde{S}(V, x_0; \theta)$ and the optimal control sequences $\tilde{p}_Q(V^*|U_B, \Sigma, x_0; \theta)$ under stochastic dynamics. We recall that under the deterministic dynamics, RHIRL uses importance sampling in equation (12) to estimate the $\frac{\partial \mathcal{L}}{\partial \theta}$ so as to update the cost function parameter θ :

$$\frac{\partial}{\partial \theta} \mathcal{L}(\theta; D, x_0) \approx \frac{1}{N} \sum_{i=1}^N \frac{1}{\lambda} \frac{\partial}{\partial \theta} S(V_i, x_0; \theta) - \frac{1}{M} \sum_{j=1}^M \frac{1}{\lambda} w(V_j) \frac{\partial}{\partial \theta} S(V_j, x_0; \theta),$$

whereas the N control sequences in the first term are from the expert demonstration D_t and the M control sequences in the second term are from our approximated optimal control sequence distribution $p_Q(V^*)$, and $w(V_j)$ is the importance sampling weight.

To estimate $\frac{\partial \mathcal{L}}{\partial \theta}$, we need to calculate/approximate the importance sampling weight $w(V)$, and the derivative of the state trajectory cost $\frac{\partial}{\partial \theta} S(V, x_0; \theta)$ w.r.t θ . We recall that the importance sampling weight $w(V)$ depends on the state trajectory cost $S(V, x_0; \theta)$ in equation (11):

$$w(V) \propto \exp\left(-\frac{1}{\lambda}\left(S(V, x_0; \theta) + \lambda \sum_{k=0}^{K-1} u_k^\top \Sigma^{-1} v_k\right)\right)$$

Under the deterministic dynamics, the importance sampling weight $w(V)$ is estimated using Monte-Carlo approximation with M state trajectory samples as follows:

$$w(V) \approx \frac{\exp\left(-\frac{1}{\lambda}(S(V, x_0; \theta) + \lambda \sum_{k=0}^{K-1} u_k^\top \Sigma^{-1} v_k)\right)}{\sum_{j=0}^{M-1} \exp\left(-\frac{1}{\lambda}(S(V_j, x_0; \theta) + \lambda \sum_{k=0}^{K-1} u_k^{j\top} \Sigma^{-1} v_k^j)\right)} \quad (56)$$

$$= \frac{\exp\left(-\frac{1}{\lambda}(\sum_{t=0}^K g(x_t; \theta) + \lambda \sum_{k=0}^{K-1} u_k^\top \Sigma^{-1} v_k)\right)}{\sum_{j=0}^{M-1} \exp\left(-\frac{1}{\lambda}(\sum_{t=0}^K g(x_t^j; \theta) + \lambda \sum_{k=0}^{K-1} u_k^{j\top} \Sigma^{-1} v_k^j)\right)} \quad (57)$$

Since the state trajectory cost $S(V, x_0; \theta)$ is a linear sum of the cost of all states, $\frac{\partial S}{\partial \theta}$ can be directly computed as follows:

$$\frac{\partial}{\partial \theta} S(V, x_0; \theta) = \frac{\partial}{\partial \theta} \sum_{t=0}^K g(x_t; \theta) = \sum_{t=0}^K \frac{\partial}{\partial \theta} g(x_t; \theta) \quad (58)$$

To extend RHIRL to stochastic dynamics, when the state trajectory cost function is now $\tilde{S}(V, x_0; \theta)$, we need to redefine how to estimate $\tilde{w}(V)$ and consequently $\frac{\partial \tilde{S}}{\partial \theta}$.

When extend to stochastic dynamic, we have the following:

$$\tilde{w}(V) \propto \exp\left(-\frac{1}{\lambda} \left(\tilde{S}(V, x_0; \theta) + \lambda \sum_{k=0}^{K-1} u_k^\top \Sigma^{-1} v_k \right)\right) \quad (59)$$

and we still adopts Monte-carlo sampling to approximate $\tilde{w}(V)$. However, since $\tilde{S}(V, x_0; \theta)$ now measures the expected state trajectory cost under the stochastic dynamics, we need to go one step further and use sampling to estimate $\tilde{S}(V, x_0; \theta)$ using M^s number of state trajectories $\tau_h = (x_0, x_1^h, \dots, x_K^h)$ per (x_0, V) pair:

$$\tilde{S}(V, x_0; \theta) \approx \frac{1}{M^s} \sum_{h=0}^{M^s-1} \sum_{t=0}^K g(x_t^h; \theta) \quad (60)$$

Therefore, now the importance sampling weight $\tilde{w}(V)$ is approximated from $M \times M^s$ state trajectories, with M^s trajectories from each (x_0, V_j) pair as follows:

$$\tilde{w}(V) \approx \frac{\exp\left(-\frac{1}{\lambda} (\tilde{S}(V, x_0; \theta) + \lambda \sum_{k=0}^{K-1} u_k^\top \Sigma^{-1} v_k)\right)}{\sum_{j=0}^{M-1} \exp\left(-\frac{1}{\lambda} (\tilde{S}(V_j, x_0; \theta) + \lambda \sum_{k=0}^{K-1} u_k^{j\top} \Sigma^{-1} v_k^j)\right)} \quad (61)$$

$$\approx \frac{\exp\left(-\frac{1}{\lambda} \left(\frac{1}{M^s} \sum_{h=0}^{M^s-1} \sum_{t=0}^K g(x_t^h; \theta) + \lambda \sum_{k=0}^{K-1} u_k^\top \Sigma^{-1} v_k\right)\right)}{\sum_{j=0}^{M-1} \exp\left(-\frac{1}{\lambda} \left(\frac{1}{M^s} \sum_{h=0}^{M^s-1} \sum_{t=0}^K g(x_t^{jh}; \theta) + \lambda \sum_{k=0}^{K-1} u_k^{j\top} \Sigma^{-1} v_k^j\right)\right)} \quad (62)$$

We emphasize that under the stochastic dynamics, we use the state trajectory samples to estimate both the state trajectory cost $\tilde{S}(V, x_0; \theta)$ and the importance sampling weight $\tilde{w}(V)$. Since the \tilde{S} now measures the expected cost over a distribution of trajectories, we need more samples to estimate $\tilde{w}(V)$ compared to the deterministic setting.

Moreover, in the final $\frac{\partial \mathcal{L}}{\partial \theta}$, we need to differentiate $\tilde{S}(V, x_0; \theta)$ w.r.t. θ . Since \tilde{S} is estimated from sampling, we have:

$$\frac{\partial}{\partial \theta} \tilde{S}(V, x_0; \theta) \approx \frac{\partial}{\partial \theta} \frac{1}{M^s} \sum_{h=0}^{M^s-1} \sum_{t=0}^K g(x_t^h; \theta) \approx \frac{1}{M^s} \sum_{h=0}^{M^s-1} \sum_{t=0}^K \frac{\partial}{\partial \theta} g(x_t^h; \theta) \quad (63)$$

Finally, we summarize how to extend RHIRL to stochastic dynamics. In stochastic dynamics, each control sequence V will map to a distribution of state trajectory $p(\tau|V, x_0)$. Hence, we adapt our measure of state trajectory cost $S(V, x_0; \theta)$ from a single trajectory to be the expected cost over a distribution of state trajectories $\tilde{S}(V, x_0; \theta)$ in equation (37). Next, we revise the optimal control sequence distribution to a stochastic setting $\tilde{p}_Q(V^*|U_B, \Sigma, x_0; \theta)$ in equation (47). More specifically, we show under the stochastic dynamics, the optimal control sequences V^* is chosen based on the expected cost of its resulting state trajectories. Finally, under this new state trajectory cost $\tilde{S}(V, x_0; \theta)$ and the optimal control sequence distribution $\tilde{p}_Q(V^*|U_B, \Sigma, x_0; \theta)$, we adapt the approximation of the importance sampling weight $\tilde{w}(V)$ and consequently the gradient of the overall loss w.r.t θ by adding one more sampling process to estimate the new state trajectory cost. Moreover, in practice, we can use the same set of samples to estimate both state trajectory cost $\tilde{S}(V, x_0; \theta)$ and the importance sampling weights $\tilde{w}(V)$.

C.4 Difference between MaxEntIRL and RHIRL in Stochastic Dynamics

To extend to stochastic dynamic, MaxEntIRL and its variants need a different measure of the entropy and redefine the overall objective function to remove the extra entropy in the original setting [35], while RHIRL simply uses more samples to approximate the state trajectory cost function. In this

section, we study why MaxEntIRL and RHIRL adapt different treatments to extend to stochastic dynamic. We first describe how MaxEntIRL extends to stochastic dynamics, then we investigate the difference between MaxEntIRL and RHIRL when extending to stochastic dynamics.

Extending MaxEntIRL to stochastic dynamic: Maximum Causal Entropy IRL [35]

MaxEntIRL resolves the fundamental ambiguity in IRL problem formulation by choosing a trajectory distribution that maximize the entropy among the distributions that matches the feature expectation of the expert. MaxEntIRL and its invariant measure the distribution in trajectory space $p(\tau)$, where $\tau = (x_0, u_0, x_1, \dots, x_T)$ is the trajectory consists of sequences of states and controls. The objective of MaxEntIRL is as follows:

$$\max H(p(\tau)) = \sum p(\tau) \ln \frac{1}{p(\tau)} \quad (64)$$

subject to the constraints

$$\mathbb{E}_{\pi^L}(\phi(\tau)) = \mathbb{E}_{\pi^E}(\phi(\tau)), \quad (65)$$

$$\sum_{\tau} p(\tau) = 1, \forall \tau, p(\tau) > 0 \quad (66)$$

whereas $\mathbb{E}_{\pi^L}(\phi(\tau))$ is the expected feature count w.r.t the learner's policy and $\mathbb{E}_{\pi^E}(\phi(\tau))$ is that of the expert. Solving this in a deterministic dynamics, the maximum entropy distribution follows

$$p(\tau) \propto \exp(R(\tau)) \quad (67)$$

whereas $p(\tau)$ is the probability of the trajectory and $R(\tau) = w^T \phi(\tau)$ is the reward of τ .

Therefore, the probability of the trajectory can be expressed as

$$p(\tau|w) = \frac{1}{Z(w)} \exp(w^T \phi(\tau)) \quad (68)$$

whereas Z is the partition function.

However, the distribution above only hold for the deterministic case. For the stochastic dynamics, the trajectory distribution is also affected by the transition probabilities, that is,

$$p(\tau|w) = \frac{1}{Z(w)} \exp(w^T \phi(\tau)) \prod_{x_{t+1}, u_t, x_t \in \tau} p(x_{t+1}|u_t, x_t) \quad (69)$$

If we follow the standard MaxEntIRL's objective, now we are optimizing

$$\tilde{R}(\tau) = w^T \phi(\tau) + \sum_t \log p(x_{t+1}|u_t, x_t) \quad (70)$$

instead of $R(\tau) = w^T \phi(\tau)$. The second term is the bias due to the stochasticity of the environment.

The implication is that this bias term will cause the learning policies to target areas in state-space of high stochasticity, instead the state-space with high reward. The fundamental problem is that the entropy is measured over the entire trajectory, therefore, how we choose the current action will be affected by the stochasticity of the future states. This is unnatural because the stochasticity in the future should not affect the current decision. Therefore, Maximum Casual Entropy IRL [35] is proposed to only account the entropy of the sub-trajectories up to the current action:

$$\begin{aligned} H(u_{1:T}|x_{1:T}) &= \sum_{t=1}^T H(u_t|u_{1:t-1}, x_{1:t}) \\ &= - \sum_{t=1}^T \sum_{u_{1:t}, x_{1:t}} p(u_{1:t}, x_{1:t}) \ln(\pi(u_t|u_{1:t-1}, x_{1:t})) \\ &= - \sum_{t=1}^T \sum_{u_{1:t}, x_{1:t}} \prod_{i=1}^t p(u_i|x_{i-1}, u_{i-1}) \ln(\pi(u_t|x_t)) \end{aligned} \quad (71)$$

Therefore, the maximum causal entropy distribution is

$$p(u_t|x_t, w) = \frac{1}{Z(w)} (\exp(w^T \phi(x_t, u_t)) - \sum_{\tau > t} \mathbb{E}_{x, u} [\log p(u_\tau|x_\tau, w)|x_t, u_t]) \quad (72)$$

Difference between Max-Entropy and RHIRL in stochastic dynamics

The key difference between MaxEntIRL and RHIRL lies in that the distribution they optimize over: MaxEntIRL maximizes the entropy in the state trajectory space, while RHIRL minimizes the "free-energy" in the control sequence space. Since the stochastic dynamic directly changes the distribution of state trajectories, when we are maximizing the entropy of the state trajectory distribution, our objective is naturally corrupted by the entropy due to the stochasticity of the dynamics. On the other hand, when we are operating on the control sequence distribution, which allows us to measure the state-trajectory cost conditioned on the given control sequence by taking an expectation over the stochastic dynamics. We can now choose the action sequence that, in expectation, results in the lowest state trajectory cost. Moreover, for each optimization step, we optimize the control sequence of length K at once, not one control by another, therefore, we do not suffer from the wrong causal effect as for MaxEntIRL.

C.5 Dependency between the Receding Horizon K and the Dynamic Model Accuracy

In this section, we study the scenario where the ground-truth dynamic f is stochastic, yet we use an approximated deterministic model \hat{f} to perform planning and reward learning. We show that, in the case when the dynamic model is inaccurate, using a receding horizon K that is shorter than the task horizon T can lead to a better planning performance. Intuitively, it shows that even if the our dynamic model used for planning lacks the capacity to accurately model the stochastic dynamics, using a receding horizon $K < T$ in practice reduces the impact of model inaccuracy on the planning performance.

For the ease of theoretical analysis, we use discount factor γ as an implicit notion of planning horizon: the larger the γ , the longer the effective planning horizon as rewards further into the future have an effect on the choice of optimal action in the current state. Therefore, instead of using the effective horizon $K < T$, we show that using a discount factor $\gamma < \gamma_{task}$ improves the planning performance.

We emphasize that in this section, we mainly investigate the planning performance loss due to the dynamic model inaccuracy. That is, given a fixed reward function R bounded by $[0, R_{max}]$, we measure the planning loss as the difference between the value function $V_{f, \gamma_{task}}$, when evaluated under the ground truth dynamic model and the discount factor, of the optimal policy computed using the ground-truth transition model and the discount factor $\pi_{f, \gamma_{task}}^*$ and that of using the estimated deterministic model and an smaller effective planning horizon $\pi_{\hat{f}, \gamma}^*$:

$$\text{planning loss: } ||V_{f, \gamma_{task}}^{\pi_{f, \gamma_{task}}^*} - V_{\hat{f}, \gamma}^{\pi_{\hat{f}, \gamma}^*}|| \quad (73)$$

Prior work [36] shows that when the transition model is estimated from data and is prone to approximation error, the policy found using a shorter planning horizon $K < T$ can actually be better than a policy learned with the true dynamics. More specifically, this paper derives a planning error bounded in the theorem below.

Theorem C.1. *Let $M = \langle S, A, R, f, \gamma_{task} \rangle$ be an MDP with non-negative rewards and evaluation discount factor γ_{task} . Let \hat{M} be an MDP comprising the transition function \hat{f} estimated from data. Then the certainty-equivalent planning with \hat{M} using effective discount factor $\gamma \leq \gamma_{task}$ has planning loss*

$$\begin{aligned} ||V_{f, \gamma_{task}}^{\pi_{f, \gamma_{task}}^*} - V_{\hat{f}, \gamma}^{\pi_{\hat{f}, \gamma}^*}|| &\leq \frac{\gamma_{task} - \gamma}{(1 - \gamma_{task})(1 - \gamma)} R_{max} \\ &\quad + \frac{2R_{max}}{(1 - \gamma)^2} \sqrt{\frac{1}{2n} \log \frac{2|S||A||\Pi_{R, \gamma}|}{\delta}} \end{aligned} \quad (74)$$

with probability at least $1 - \delta$.

In the equation above, $|\Pi_{R,\gamma}|$ is a complexity measure of the policy class under the fixed reward R and effective discount factor γ . Formally, $\Pi_\gamma = \{\pi : \exists R \text{ s.t. } \pi \text{ is optimal in } \langle S, A, T, R, \gamma \rangle\}$, hence $|\Pi_{R,\gamma}|$ measures the number of the potentially optimal policies given a certain γ . The paper [36] proved that $|\Pi_{R,\gamma}|$ increases dramatically as γ increases.

In our scenario, using a deterministic model \hat{f} to approximate the stochastic dynamic f results in model inaccuracy, therefore, we can use the error bound in Theorem C.1. The planning error bound in Theorem C.1 has two terms: the first term measures the planning error due to using the effective discount factor $\gamma < \gamma_{task}$, and the second term measures the planning loss due to planning using the inaccurate dynamic model \hat{f} . When γ increases and approaches γ_{task} , the first error term decreases while the second error term increases. Therefore, since the two error terms move in different direction as γ increases from 0 to γ_{task} , the optimal effective horizon γ^* that minimizes the planning loss in equation (74) must be smaller than γ_{task} . That is, we define

$$\gamma^* = \arg \min_{\gamma \in [0, \gamma_{task}]} \|V_{f, \gamma_{task}}^{\pi_{\hat{f}, \gamma_{task}}} - V_{\hat{f}, \gamma}^{\pi_{\hat{f}, \gamma}}\| \quad (75)$$

then, $\gamma^* < \gamma_{task}$. Therefore, even if we directly apply the standard RHIRL in a stochastic dynamics, our receding horizon controller that plans with an effective horizon $K \ll T$ practically reduces the extent of the planning error due to model inaccuracy. In general, when the dynamic model is inaccurate, which is the case in most real-life problems, RHIRL is inherently better than other IRL algorithms, as it has the flexibility to control the planning horizon so as to reduce the impact of the model inaccuracy on the performance of policy optimization.

D Experimental Details

In this section, we list down the implementation details of RHIRL and the baselines. We also report the hyperparameters used in the experiments, the detailed network architectures, training procedures and evaluation procedures used for our experiments.

D.1 Practical Issues of RHIRL

Control Noise Covariance Approximation

The actual control noise covariance Σ is unknown to RHIRL and the baselines. However, RHIRL uses the noise covariance matrix Σ to sample the controls around the nominal control (Algorithm 1, line 6) and calculate the quadratic control cost in Equation (37). Since we have no access to the true Σ , RHIRL approximates Σ as a constant factor of the identity matrix βI , whereas β is the hyperparameter we optimize using grid search and I is the identity matrix with its width equals to the dimension of the action space. Therefore, instead of sampling the controls from $\mathcal{N}(V|U, \Sigma)$, we sample from $\mathcal{N}(V|U, \beta I)$ in practice. We also use βI in Equation (37) to replace the unknown Σ . Even in the noise-free environment, we set β to a non-zero value to foster exploration; otherwise, the importance sampling degenerates to the single nominal control.

Our experiment shows that RHIRL is robust to the choice of β : the cost learning performs well even if $\beta I \neq \Sigma$. This may be attributed to the fact that we jointly optimize the state cost function and β . Therefore the learned state cost function may compensate for the inaccurate approximation for Σ .

Numerical Stability

Equation(11) forms the basis of importance sampling and estimation. However, the learned cost can be a huge negative number, which causes numerical instability in estimating the importance weights. To mitigate this issue, we subtract the minimum trajectory cost S_{min} from all rollouts to improve the numeric stability. Since subtracting the same number from all rollouts does not change the order of the preference, this operation does not affect the optimality of our derivation.

Nominal Control Initialization and Local Optimality

RHIRL samples around the nominal control sequence to collect the samples for importance sampling. However, if the initial nominal control sequence performs poorly, it is not easy to generate any good samples to improve the current control sequences. To mitigate this problem, we add an exploration strategy to the sampling process in (Algorithm 1, line 6): with probability α , we continue with the

standard sampling strategy to sample around the nominal control; with probability $1 - \alpha$ we sample uniformly from the entire action space. This helps RHIRL to correct from the unsuitable nominal control initialization and also helps RHIRL to escape the local optimal solution. We set $\alpha = 0.5$ for all tasks in our experiments.

Control Smoothness

Updating the optimal control by importance sampling might cause some jerk in the control space. In order to make the control change smoothly in its local space, we apply a Savitzky–Golay filter over the time horizon dimension to constrain the control that does not change too much over the time horizon.

D.2 Training Details

We list the hyper-parameters of RHIRL for different tasks. These hyper-parameters were selected via grid search.

| Task | K | β | batch size | λ | lr | weight decay |
|--------------------------|----|---------|------------|-----------|------|--------------|
| Pendulum-v0 | 20 | 0.8 | 50 | 0.10 | 1e-4 | 8e-5 |
| LunarLanderContinuous-v2 | 40 | 0.6 | 200 | 0.10 | 1e-4 | 8e-5 |
| Hopper-v2 | 20 | 0.8 | 100 | 0.10 | 1e-4 | 8e-5 |
| Walker2d-v2. | 30 | 0.6 | 150 | 0.10 | 1e-4 | 8e-5 |
| Ant-v2 | 15 | 1.2 | 200 | 0.10 | 1e-4 | 8e-5 |
| CarRacing-v0 | 15 | 1.0 | 200 | 0.10 | 1e-4 | 8e-5 |

The implementation of the baselines (f-IRL, AIRL and GAIL) are adapted from f-IRL’s [20] official repository. We use the hyperparameters reported in f-IRL for the MuJoCo tasks and performed grid search on the hyperparameters for the rest of the tasks. SAC[31] is used as the base MaxEnt RL algorithm for both expert policy and the baselines optimization algorithm. We use a tanh squashed Gaussian as the policy network for Pendulum-v0, LunarLander-v2, and the MuJoCo tasks; and we use a Gaussian Convolutional policy as the policy network for CarRacing-v0. The mean and std of the Gaussian are parameterized by a ReLU MLP of size (64, 64). Adam is used as the optimizer. We use the reported SAC temperature, $\alpha = 0.2$, reward scale $c = 0.2$, and gradient penalty coefficient $\lambda = 4.0$. The rest of the hyperparameters for f-IRL, GAIL and AIRL are listed below.

| Task | SAC learning rate | SAC replay buffer size | Reward/Value model learning rate | l2 weight decay |
|--------------------------|-------------------|------------------------|----------------------------------|-----------------|
| Pendulum-v0 | 1e-4 | 100000 | 1e-5 | 1e-3 |
| LunarlanderContinuous-v2 | 1e-3 | 100000 | 1e-5 | 1e-3 |
| Hopper-v2 | 1e-5 | 1000000 | 1e-5 | 1e-3 |
| Walker2d-v2 | 1e-5 | 1000000 | 1e-5 | 1e-3 |
| Ant-v2 | 3e-4 | 1000000 | 1e-4 | 1e-3 |
| CarRacing-v0 | 4e-4 | 10000000 | 1e-4 | 1e-3 |

D.3 Reward Function/Discriminators Network architectures

We use the same neural network architecture to parameterize the cost-function/reward-function/discriminator for all methods. For continuous control task with raw state input, i.e. pendulum, lunarlander and the MuJoCo tasks, we use two-layer of MLP with ReLU activation function to parameterized the cost function/discriminator. The hidden size for Pendulum-v0 is (32, 32), and (64, 64) for the rest of the tasks.

For continuous control task with image input, i.e. carracing, we use a four convolutional layer with kernel size 3×3 as the feature extractor. The output of the CNN layer is vector with size (128,) and is fed into the same reward network as describe above.

Table 3: Performance of RHIRL, f-IRL, GAIL, and AIRL. We report the mean and the standard deviation of the policy returns using the ground-truth task reward. Higher values indicate better performance.

| | | No Noise $\Sigma = 0$ | Mild Noise $\Sigma = 0.2$ | High Noise $\Sigma = 0.5$ |
|-------------|--------|--------------------------|------------------------------|------------------------------|
| Pendulum | Expert | -154.69 \pm 50.05 | -156.50 \pm 70.72 | -168.54 \pm 80.89 |
| | RHIRL | -125.95 \pm 1.21 | -122.33 \pm 3.44 | -132.39 \pm 10.36 |
| | f-IRL | -121.94 \pm 97.21 | -127.51 \pm 104.55 | -197.36 \pm 106.92 |
| | AIRL | -131.64 \pm 1.16 | -184.62 \pm 88.16 | -203.12 \pm 80.57 |
| | GAIL | -207.05 \pm 57.41 | -207.14 \pm 57.52 | -253.85 \pm 181.84 |
| LunarLander | Expert | 235.13 \pm 43.59 | 222.65 \pm 56.35 | 164.52 \pm 36.79 |
| | RHIRL | 246.39 \pm 10.96 | 233.73 \pm 23.75 | 198.23 \pm 47.8 |
| | f-IRL | 179.03 \pm 9.19 | 141.73 \pm 11.81 | 121.67 \pm 22.77 |
| | AIRL | 174.49 \pm 35.17 | 132.76 \pm 85.59 | 95.61 \pm 19.25 |
| | GAIL | 169.98 \pm 15.43 | 125.5 \pm 16.78 | 100.24 \pm 79.04 |
| Hopper | Expert | 3222.48 \pm 390.65 | 3159.32 \pm 520.00 | 2887.72 \pm 483.93 |
| | RHIRL | 3071.63 \pm 122.03 | 3121.72 \pm 278.98 | 2776.2 \pm 345.90 |
| | f-IRL | 3080.34 \pm 458.96 | 2580.19 \pm 637.21 | 1270.24 \pm 539.84 |
| | AIRL | 18.9 \pm 0.79 | 33.52 \pm 3.86 | 18.38 \pm 7.84 |
| | GAIL | 2642.59 \pm 187.33 | 1576.25 \pm 1051.98 | 702.33 \pm 151.37 |
| Walker2d | Expert | 4999.47 \pm 55.99 | 4500.43 \pm 114.48 | 3624.48 \pm 95.05 |
| | RHIRL | 4939.44 \pm 100.28 | 4473.332 \pm 324.34 | 3446.55 \pm 507.89 |
| | f-IRL | 4927.92 \pm 529.95 | 3697.36 \pm 711.56 | 2831.91 \pm 993.76 |
| | AIRL | -2.51 \pm 0.69 | 22.24 \pm 10.74 | 6.5 \pm 5.03 |
| | GAIL | 2489.04 \pm 813.31 | 2884.35 \pm 59.88 | 1840.62 \pm 778.3 |
| Ant | Expert | 5759.22 \pm 173.57 | 2557.37 \pm 501.95 | 252.62 \pm 91.44 |
| | RHIRL | 4987.67 \pm 149.2 | 2373.32 \pm 529.3 | 230.8 \pm 253.39 |
| | f-IRL | 5022.42 \pm 108.07 | 2034.87 \pm 262.29 | 197.2 \pm 200.45 |
| | AIRL | 1000.4 \pm 0.79 | 849.05 \pm 30.15 | -7.43 \pm 6.01 |
| | GAIL | 2784.87 \pm 301.66 | 1022.04 \pm 580.49 | -416.69 \pm 292.23 |
| CarRacing | Expert | 903.25 \pm 0.23 | 702.01 \pm 0.3 | 281.12 \pm 0.34 |
| | RHIRL | 359.61 \pm 40.32 | 206.21 \pm 19.87 | 53.97 \pm 3.24 |
| | f-IRL | 85.45 \pm 47.4 | 18.32 \pm 27.89 | 2.04 \pm 13.8 |
| | AIRL | -21.97 \pm 2.67 | -25.25 \pm 5.98 | -32.31 \pm 7.43 |
| | GAIL | 2.62 \pm 3.41 | -7.65 \pm 4.77 | -15.88 \pm 5.89 |

D.4 Numerical Experiments Results

We report the average returns and the standard deviation for Table 1 and Table 2 in Table 3 and Table 4 respectively. The mean and standard deviation computed from 3 trials for each entry of the tables.

Table 4: Generalization of learned cost functions over different noise levels.

| | | Noise-free for learning | Noise Level Σ for Testing | |
|-------------|-------|-------------------------|----------------------------------|-----------------------|
| | | | 0.2 | 0.5 |
| Pendulum | RHIRL | -125.95 \pm 1.21 | -125.01 \pm 4.53 | -126.4 \pm 7.73 |
| | f-IRL | -121.94 \pm 97.21 | -199.44 \pm 96.99 | -220.74 \pm 79.75 |
| | AIRL | -131.64 \pm 1.16 | -247.86 \pm 11.44 | -304.48 \pm 20.78 |
| | GAIL | -207.05 \pm 57.41 | -220.6 \pm 69.82 | -270.81 \pm 79.68 |
| LunarLander | RHIRL | 246.39 \pm 10.96 | 205.66 \pm 24.67 | 175.82 \pm 52.12 |
| | f-IRL | 179.03 \pm 9.19 | 121.80 \pm 20.94 | 102.06 \pm 22.31 |
| | AIRL | 174.49 \pm 35.17 | 31.46 \pm 9.68 | 22.29 \pm 14.01 |
| | GAIL | 169.98 \pm 15.43 | 101.80 \pm 23.12 | 78.33 \pm 24.15 |
| Hopper | RHIRL | 3071.63 \pm 122.03 | 2577.28 \pm 409.33 | 2152.08 \pm 342.21 |
| | f-IRL | 3080.34 \pm 458.96 | 2110.52 \pm 26.71 | 1984.29 \pm 31.88 |
| | AIRL | 18.9 \pm 0.79 | 18.86 \pm 4.80 | 8.78 \pm 10.89 |
| | GAIL | 2642.59 \pm 187.33 | 215.29 \pm 27.76 | 132.15 \pm 30.20 |
| Walker2d | RHIRL | 4939.44 \pm 100.28 | 4039.44 \pm 39.2 | 3440.23 \pm 531.08 |
| | f-IRL | 4927.92 \pm 529.95 | 2976.66 \pm 396.57 | 1090.11 \pm 1389.56 |
| | AIRL | -2.51 \pm 0.69 | 1380.84 \pm 364.95 | 1787.15 \pm 230.94 |
| | GAIL | 2489.04 \pm 813.31 | 103.15 \pm 121.84 | 124.15 \pm 82.84 |
| Ant | RHIRL | 4987.67 \pm 149.2 | 3192.82 \pm 162.12 | 867.08 \pm 204.28 |
| | f-IRL | 5022.42 \pm 108.07 | 2042.41 \pm 129.89 | 472.77 \pm 110.2 |
| | AIRL | 1000.4 \pm 0.79 | 845.69 \pm 29.01 | 0.69 \pm 20.49 |
| | GAIL | 2784.87 \pm 301.66 | -6.41 \pm 21.17 | -79.89 \pm 142.43 |
| CarRacing | RHIRL | 359.61 \pm 40.32 | 261.78 \pm 54.44 | 110.12 \pm 58.90 |
| | f-IRL | 85.45 \pm 47.4 | 16.12 \pm 67.82 | -24.78 \pm 2.12 |
| | AIRL | -21.97 \pm 2.67 | -27.09 \pm 6.65 | -23.96 \pm 4.11 |
| | GAIL | 2.62 \pm 3.41 | -6.41 \pm 3.22 | -49.89 \pm 7.98 |