CROSS-DOMAIN POLICY OPTIMIZATION VIA BELLMAN CONSISTENCY AND HYBRID CRITICS

Anonymous authorsPaper under double-blind review

000

001

002 003 004

010 011

012

013

014

016

017

018

019

021

025

026

027

028029030

031

033

034

035

037

040

041

042

043

044

046

047

048

051

052

ABSTRACT

Cross-domain reinforcement learning (CDRL) is meant to improve the data efficiency of RL by leveraging the data samples collected from a source domain to facilitate the learning in a similar target domain. Despite its potential, cross-domain transfer in RL is known to have two fundamental and intertwined challenges: (i) The source and target domains can have distinct state space or action space, and this makes direct transfer infeasible and thereby requires more sophisticated interdomain mappings; (ii) The transferability of a source-domain model in RL is not easily identifiable a priori, and hence CDRL can be prone to negative effect during transfer. In this paper, we propose to jointly tackle these two challenges through the lens of cross-domain Bellman consistency and hybrid critic. Specifically, we first introduce the notion of cross-domain Bellman consistency as a way to measure transferability of a source-domain model. Then, we propose QAvatar, which combines the Q functions from both the source and target domains with an adaptive hyperparameter-free weight function. Through this design, we characterize the convergence behavior of QAvatar and show that QAvatar achieves reliable transfer in the sense that it effectively leverages a source-domain Q function for knowledge transfer to the target domain. Through experiments, we demonstrate that QAvatar achieves favorable transferability across various RL benchmark tasks, including locomotion and robot arm manipulation.

1 Introduction

Cross-domain reinforcement learning (CDRL) serves as a practical framework to improve the sample efficiency of RL from the perspective of transfer learning, which leverages the pre-trained models from a source domain to enable knowledge transfer to the target domain, under the presumption that the data collection and model training are much less costly in the source domain (e.g., simulators). A plethora of the existing CDRL methods focuses on knowledge transfer across environments that share the same state-action spaces but with different transition dynamics. This setting has been extensively studied from a variety of perspectives, such as reward augmentation (Eysenbach et al., 2021; Liu et al., 2022), data filtering (Xu et al., 2023), and latent representations (Lyu et al., 2024). Despite the above progress, to fully realize the promise of CDRL, there are two fundamental challenges to tackle: (i) Distinct state and/or action spaces between domains: To support flexible transfer across a wide variety of domains, the generic CDRL is required to address the discrepancies in the state and action spaces between source and target domains. Take robot control as an example. One common scenario is to apply direct policy transfer between robot agents of different morphologies (Zhang et al., 2021), which naturally leads to a discrepancy in representations. This discrepancy significantly complicates the transfer of either data samples or learned source-domain models. (ii) Unknown transferability of a source-domain model to the target domain: CDRL conventionally presumes that the source-domain model can achieve effective transfer under a properly learned cross-domain correspondence. However, in practice, given that the data budget of the target domain is limited, it is rather difficult to determine a priori the transferability of a source-domain model. Indeed, it has been widely observed that transfer learning from the source domain can have a negative impact on the target domain (Weiss et al., 2016; Pan & Yang, 2009).

As a consequence, despite that CDRL has been shown to succeed in various scenarios, without a proper design, the performance of CDRL could actually be much worse than the vanilla target-domain model learned without using any source knowledge. Notably, to tackle (i), several approaches have

been proposed to address such representation discrepancy by learning state-action correspondence, either in the typical RL (You et al., 2022) or unsupervised settings (Zhang et al., 2021; Gui et al., 2023). However, existing solutions are all oblivious to the issues of model transferability between the domains. Hence, one fundamental research question about CDRL remains largely open:

How to achieve effective transfer in CDRL under distinct state-action spaces without the knowledge of the transferability of the pre-trained source-domain model?

In this paper, we affirmatively address the above question by revisiting cross-domain state-action correspondence through the lens of *cross-domain Bellman consistency*, which quantifies the transferability of a source-domain model. To enable reliable transfer across varying levels of source-model transferability, we introduce a novel CDRL framework, *QAvatar*, which integrates source-domain and target-domain critics. Drawing an analogy from the movie *Avatar*, where humans remotely control genetically engineered bodies to adapt to alien environments, *Q*Avatar updates the target-domain policy via a weighted combination of the target- and source-domain Q functions, while learning the state-action correspondence by minimizing a cross-domain Bellman loss.

To validate this idea, we first present a tabular prototype of QAvatar and show that it attains a tight sub-optimality bound under an adaptive, hyperparameter-free weight function, regardless of source model transferability. This ensures improved sample efficiency while avoiding poor transfer. Building on this, we develop a practical version by combining QAvatar with a normalizing flow-based mapping for learning state-action correspondence.

The main contributions of this paper can be summarized as follows: 1) We propose the QAvatar framework that achieves knowledge transfer between two domains with distinct state and action spaces for improving sample efficiency. We then present a prototypical QAvatar algorithm and establish its convergence property. 2) We further substantiate the QAvatar framework by proposing a practical implementation with a normalizing-flow-based state-action mapping. This further demonstrates the compatibility of QAvatar with off-the-shelf methods for learning state-action correspondence. 3) Through experiments and an ablation study, we show that QAvatar outperforms the CDRL benchmark algorithms on various RL benchmark tasks.

2 RELATED WORK

CDRL across domains with distinct state and action spaces. The existing approaches can be divided into two main categories: (i) Manually designed latent mapping: In (Ammar & Taylor, 2012; Gupta et al., 2017; Ammar et al., 2012), the trajectories are mapped manually from the source domain and the target domain to a common latent space. The distance between latent states can then be calculated to find the correspondence of the states from the different domains. (ii) Learned inter-domain mapping: In (Taylor et al., 2008; Zhang et al., 2021; You et al., 2022; Gui et al., 2023; Zhu et al., 2024), the inter-domain mapping is mainly learned by enforcing dynamics alignment (or termed dynamics cycle consistency in (Zhang et al., 2021)). Additional properties have also been incorporated as auxiliary loss functions in learning the inter-domain mapping, including domain cycle consistency (Zhang et al., 2021), effect cycle consistency (Zhu et al., 2024), maximizing mutual information between states and embeddings (You et al., 2022) However, the existing approaches all presume that the domains are sufficiently similar and do not have any performance guarantees. By contrast, we propose a reliable CDRL method that can achieve transfer regardless of source-domain model quality or domain similarity with guarantees.

CDRL across domains with *identical* state and action spaces. Various methods have been proposed for the case where source and target domains share the same state and action spaces but are subject to dynamics mismatch. Existing methods include (i) using the samples from both source and target domains jointly for learning (Eysenbach et al., 2021; Liu et al., 2022; Xu et al., 2023; Lyu et al., 2024), (ii) explicit characterization of domain similarity (Behboudian et al., 2022; Sreenivasan et al., 2023), and (iii) using both Q-functions for Q-learning updates (Wang et al., 2020). However, given the assumption on identical state-action spaces, they are not readily applicable to our CDRL setting.

3 PRELIMINARIES

In this section, we provide the problem statement and basic building blocks of CDRL as well as the useful notation needed by subsequent sections. For a set \mathcal{X} , we let $\Delta(\mathcal{X})$ denote the set of probability distributions over \mathcal{X} . As in typical RL, we model each environment as an infinite-horizon

discounted Markov decision process (MDP) denoted by $\mathcal{M}:=(\mathcal{S},\mathcal{A},P,r,\gamma,\mu)$, where (i) \mathcal{S} and \mathcal{A} represent the state space and action space, (ii) $P:\mathcal{S}\times\mathcal{A}\to\Delta(\mathcal{S})$ denotes the transition function, (iii) $r:\mathcal{S}\times\mathcal{A}\to[0,1]$ is the reward function (without loss of generality, we presume the rewards lie in the [0,1] interval), (iv) $\gamma\in[0,1)$ is the discounted factor, and (v) $\mu\in\Delta(\mathcal{S}\times\mathcal{A})$ denotes the initial state-action distribution. Notably, the use of an initial distribution over states and actions is a standard setting in the literature of natural policy gradient (NPG) (Agarwal et al., 2021a; Ding et al., 2020; Yuan et al., 2022; Agarwal et al., 2020; Zhou et al., 2024). Given any policy $\pi:\mathcal{S}\to\Delta(\mathcal{A})$, let $\tau=(s_0,a_0,r_1,\cdots)$ denote a (random) trajectory generated under π in \mathcal{M} , and the expected total discounted reward under π is $V_{\mathcal{M}}^{\pi}(\mu):=\mathbb{E}[\sum_{t=0}^{\infty}\gamma^t r(s_t,a_t)|\pi;s_0,a_0\sim\mu]$. We use $Q_{\mathcal{M}}^{\pi}(s,a)$ and $V_{\mathcal{M}}^{\pi}(s)$ to denote the Q function and value function of a policy π . We also define the stateaction visitation distribution (also known as the occupancy measure in the MDP literature) of π as $d^{\pi}(s,a):=(1-\gamma)(\mu(s,a)+\sum_{t=1}^{\infty}\gamma^t\mathbb{P}(s_t=s,a_t=a;\pi,\mu))$, for each (s,a).

Problem Statement of Cross-Domain RL. In typical CDRL, the knowledge transfer involves two MDPs, namely the source-domain MDP $\mathcal{M}_{src}:=(\mathcal{S}_{src},\mathcal{A}_{src},P_{src},r_{src},\gamma,\mu_{src})$ and the target-domain MDP $\mathcal{M}_{tar}:=(\mathcal{S}_{tar},\mathcal{A}_{tar},P_{tar},r_{tar},\gamma,\mu_{tar})^1$. Notably, in addition to distinct state and action spaces, the two domains can have different reward functions, transition dynamics, and initial distributions. We assume that the two MDPs share the same discounted factor γ , which is rather mild. Moreover, the trajectories of the two domains are completely unpaired. Let Π_{tar} be the set of all stationary Markov policies for \mathcal{M}_{tar} .

The goal of the RL agent is to learn a policy π^* in the target domain such that the expected total discounted reward is maximized, *i.e.*, $\pi^* := \arg\max_{\pi \in \Pi_{tar}} V_{\mathcal{M}_{tar}}^{\pi}(\mu_{tar})$. To improve sample efficiency via knowledge transfer (compared to learning from scratch), in CDRL, the target-domain agent is granted access to $(\pi_{src}, Q_{src}, V_{src})$, which denotes a policy and the corresponding Q and value functions pre-trained in \mathcal{M}_{src} . Notably, we make no assumption on the quality of π_{src} (and hence π_{src} may not be optimal to \mathcal{M}_{src}), despite that π_{src} shall exhibit acceptable performance in practice.

In this paper, we focus on designing a reliable CDRL algorithm in that it effectively leverages a source-domain Q function $Q_{\rm src}$ for knowledge transfer to the target domain, regardless of the quality of $Q_{\rm src}$ and domain similarity.

Inter-Domain Mapping Functions. To address the discrepancy in state-action spaces in CDRL, learning an inter-domain mapping is one common block of many CDRL algorithms. Specifically, there are a variety of ways to construct the mapping functions, such as handcrafted functions (Ammar & Taylor, 2012), encoders and decoders trained by cycle consistency You et al. (2022) like cycle-GAN (Zhu et al., 2017), neural networks trained by dynamics alignment of the MDPs (Gui et al., 2023). Moreover, mapping functions have various candidate target spaces, such as a latent space, state or action spaces of the target domain (*i.e.*, from $S_{\rm src}$, $A_{\rm src}$ to $S_{\rm tar}$, $A_{\rm tar}$), and state or action spaces of the source domain (*i.e.*, from $S_{\rm tar}$, $A_{\rm tar}$ to $S_{\rm src}$, $A_{\rm src}$).

For example, Gui et al. (2023) proposed learning two mappings, $G_1: \mathcal{S}_{tar} \to \mathcal{S}_{src}$ and $G_2: \mathcal{A}_{src} \to \mathcal{A}_{tar}$, via dynamics alignment, which infers the unknown mapping between unpaired trajectories of \mathcal{M}_{src} and \mathcal{M}_{tar} by aligning one-step state transitions. However, this unsupervised approach provides no performance guarantee and can suffer from identification issues. By contrast, we propose learning inter-domain state and action mappings, $\phi: \mathcal{S}_{tar} \to \mathcal{S}_{src}$ and $\psi: \mathcal{A}_{tar} \to \mathcal{A}_{src}$, using a cross-domain Bellman-like loss with guarantees (Section 4). Appendix D.1 shows a toy example where cycle consistency fails, but the Bellman-like loss leverages target rewards to learn a better mapping.

Tabular Approximate Q-Natural Policy Gradient. Natural Policy Gradient (NPG) (Kakade, 2001; Agarwal et al., 2019) is a classical RL algorithm. In this paper, we adopt NPG under two assumptions to analyze CDRL: (i) **Tabular setting:** finite state and action spaces, with independent parameters for each state-action pair (s, a); (ii) **Approximate Q-function:** the true Q^{π} is inaccessible due to limited data, so we use an empirical approximation from samples. At iteration t, we first collect data $\mathcal{D}^{(t)}$ by executing $\pi^{(t)}$, then obtain $Q^{(t)}$ by minimizing the standard TD loss for least-squares policy evaluation (LSPE) (Lagoudakis & Parr, 2001; Yu & Bertsekas, 2009; Lazaric et al., 2012)²

$$\mathcal{L}_{\text{TD}}(Q^{(t)}; \pi^{(t)}, \mathcal{D}^{(t)}) := \hat{\mathbb{E}}_{(s, a, r, s') \in \mathcal{D}^{(t)}} \Big[\Big| r + \gamma \mathbb{E}_{a' \sim \pi^{(t)}} [Q^{(t)}(s', a')] - Q^{(t)}(s, a) \Big|^2 \Big]. \tag{1}$$

¹Throughout this paper, we use the subscripts "src" and "tar" to represent the objects in the source and target domains, respectively.

²LSPE under linear function approximation includes the tabular case via one-hot features:

Finally, we perform a one-step policy improvement: $\pi^{(t+1)}(a|s) \propto \pi^{(t)}(a|s) \exp \left(\eta Q^{(t)}(s,a)\right)$, where η is the learning rate. This update improves the policy while staying close to the original.

Notation. Throughout this paper, for any policy π and any real-valued function $h: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$, we use $h(s,\pi)$ and $\bar{h}(s,a;\pi)$ as the shorthand for $\mathbb{E}_{a \sim \pi(\cdot|s)}[h(s,a)]$ and $h(s,a) - \mathbb{E}_{a \sim \pi(\cdot|s)}[h(s,a)]$, respectively. For any real vector z and $p \ge 1$, we let $||z||_p$ be the ℓ_p -norm of z. For any real-valued function $f: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$, we use $\|f\|_{d^{\pi(t)}}$ as the shorthand for $\mathbb{E}_{(s,a) \sim d^{\pi(t)}} \left[f(s,a) \right]$.

METHODOLOGY

162

163

164

166

167

168 169 170

171

172

173

174

175

176

177

178

179

181

182

183

184

185 186

187

188

189

190

191

192 193

194 195 196

197

199 200 201

202 203 204

205 206

207

208

209

210

211

212

213

214

215

In this section, we first describe the concept of cross-domain Bellman consistency and accordingly propose the QAvatar framework in the tabular setting (i.e., S_{tar} and A_{tar} are finite). We then extend this framework to a practical deep RL implementation.

4.1 Cross-Domain BELLMAN CONSISTENCY

Algorithm 1 Direct Q Transfer (DQT)

Require: Source-domain Q function $Q_{\rm src}$, total iterations T, and $\eta = (1 - \gamma)\sqrt{1/T}$.

- 1: Initialize $\pi^{(1)}$ as a uniformly random policy.
- 2: for iteration $t=1,\cdots,T$ do 3: Select $\phi^{(t)}$ and $\psi^{(t)}$
- Update target-domain policy as in (3).
- 5: end for
- 6: **Return** $\pi_{\text{tar}}^{(T)} \sim \text{Uniform}(\{\pi^{(1)}, \cdots, \pi^{(T)}\}).$

To motivate Source domain Q-function transfer, we present the sub-optimal gap of traditional NPG. First, we describe the definitions of state-action distribution coverage and TD error.

Definition 1 (Coverage). Given a target-domain policy π^{\dagger} in \mathcal{M}_{tar} , we say that π^{\dagger} has coverage $C_{\pi^{\dagger}}$ if for any policy $\pi \in \Pi_{tar}$, we have $\|d^{\pi^{\dagger}}/d^{\pi}\|_{\infty} \leq C_{\pi^{\dagger}}$.

Assumption 1. The initial distribution is exploratory, i.e., $\mu_{tar}(s, a) > 0$, for all s, a.

Notably, $C_{\pi^{\dagger}}$ is finite if $\|d^{\pi^{\dagger}}/\mu_{\text{tar}}\|_{\infty}$ is finite (since $\|\mu_{\text{tar}}/d^{\pi}\|_{\infty} \leq 1/(1-\gamma)$ for all π by the definition of d^{π}), which holds under an exploratory initial distribution with $\mu_{tar}(s,a)>0$ for all (s, a)—a standard assumption in the NPG literature (Agarwal et al., 2021a; Ding et al., 2020; Yuan et al., 2022; Agarwal et al., 2020; Zhou et al., 2024). Intuitively, coverage enables direct comparison of Bellman errors between policies. We also use $\mu_{\text{tar,min}}$ as shorthand for $\min_{s,a} \mu_{\text{tar}}(s,a)$.

Definition 2 (TD Error). For each state-action pair (s,a) and $t \in \mathbb{N}$, the TD error $\epsilon_{td}^{(t)}(s,a)$ is $\textit{defined as } \epsilon_{\textit{td}}^{(t)}(s,a) := \big|Q_{\textit{tar}}^{(t)}(s,a) - r_{\textit{tar}}(s,a) - \gamma \mathbb{E}_{s' \sim P_{\textit{tar}}(\cdot|s,a),a' \sim \pi^{(t)}(\cdot|s')}[Q_{\textit{tar}}^{(t)}(s',a')]\big|.$

Proposition 1. Under the tabular and approximate-Q settings, and Assumption 1, the average sub-optimality of Q-NPG over T iterations is upper bounded by

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{s \sim \mu_{tar}} \left[V^{\pi^*}(s) - V^{\pi^{(t)}}(s) \right] \\
\leq \underbrace{\frac{\left[\log |\mathcal{A}_{tar}| + 1 \right]}{\sqrt{T}(1-\gamma)}}_{(a)} + \underbrace{\frac{C_0}{T} \sum_{t=1}^{T} \left\| \left| Q_{tar}^{(t)} - Q^{\pi^{(t)}} \right| \right\|_{d^{\pi^{(t)}}}}_{(b)} \leq \underbrace{\frac{\left[\log |\mathcal{A}_{tar}| + 1 \right]}{\sqrt{T}(1-\gamma)}}_{(a)} + \underbrace{\frac{C_1}{T} \sum_{t=1}^{T} \left\| \epsilon_{td}^{(t)} \right\|_{d^{\pi^{(t)}}}}_{(c)}, \quad (2)$$

where $C_0 := 2C_{\pi^*}/(1-\gamma)$ and $C_1 := 2C_{\pi^*}/((1-\gamma)^3\mu_{tar,n})$

The detailed proof of Proposition 1 is provided in Appendix B. The upper bound of the sub-optimality gap has two parts. Term (a) characterizes Q-NPG learning and converges at $O(1/\sqrt{T})$, while term (b) (or equivalently term (c)) accounts for approximation error at each iteration, which can be made arbitrarily small with enough samples (Agarwal et al., 2021a). In CDRL, limited data amplifies term (b), potentially preventing convergence to the optimal policy. To mitigate this issue, instead of learning $Q^{(t)}$ from scratch to approximate $Q^{\pi^{(t)}}$, we leverage a pre-trained source-domain Qfunction $Q_{src}(\phi^{(t)}(s), \psi^{(t)}(a))$ with inter-domain mapping $\phi^{(t)}$ and $\psi^{(t)}$ to approximate $Q^{\pi^{(t)}}$. Here, the inter-domain mappings $\phi^{(t)}$ and $\psi^{(t)}$ are introduced to address the state–action representation mismatch. For more specifically, we present Direct Q Transfer (DQT) method, in each iteration t, DQT proceeds in two steps: (i) It first updates $\phi^{(t)}$ and $\psi^{(t)}$, e.g., by gradient descent on some loss

function. (ii) The policy is updated by an NPG policy improvement step based on the pre-trained source-domain $Q_{\rm src}$ and inter-domain mappings $\phi^{(t)}, \psi^{(t)}$ as

$$\pi^{(t+1)}(a|s) \propto \pi^{(t)}(a|s) \exp\left(\eta Q_{\rm src}(\phi^{(t)}(s), \psi^{(t)}(a))\right),$$
 (3)

where η is the step size. The pseudo code is in Algorithm 1. Before characterizing the convergence behavior, we describe the cross-domain Bellman error used in Proposition 2.

Definition 3 (Cross-Domain Bellman Error). Given a pre-trained source-domain Q_{src} , inter-domain correspondences ϕ, ψ , and target-domain policy π , for each state-action pair (s,a), the cross-domain Bellman error is defined as $\epsilon_{cd}(s,a;\phi,\psi,Q_{src},\pi) := |Q_{src}(\phi(s),\psi(a)) - r_{tar}(s,a) - \gamma \mathbb{E}_{s' \sim P_{tar}(\cdot|s,a),a' \sim \pi(\cdot|s')}[Q_{src}(\phi(s'),\psi(a'))]|$.

Proposition 2. Under the DQT method in Algorithm 1 and Assumption 1, the average sub-optimality over T iterations is upper bounded as

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{s \sim \mu_{tar}} \left[V^{\pi^*}(s) - V^{\pi^{(t)}}(s) \right] \leq \underbrace{\frac{\left[\log |\mathcal{A}_{tar}| + 1 \right]}{\sqrt{T}(1 - \gamma)}}_{(a)} + \underbrace{\frac{C_0}{T} \sum_{t=1}^{T} \left\| \left| Q_{src}(\phi^{(t)}, \psi^{(t)}) - Q^{\pi^{(t)}} \right| \right\|_{d^{\pi^{(t)}}}}_{(b)} \\
\leq \underbrace{\frac{\left[\log |\mathcal{A}_{tar}| + 1 \right]}{\sqrt{T}(1 - \gamma)}}_{(a)} + \underbrace{\frac{C_1}{T} \sum_{t=1}^{T} \left\| \epsilon_{cd}(Q_{src}, \phi^{(t)}, \psi^{(t)}) \right\|_{d^{\pi^{(t)}}}}_{(c)}}_{(c)}, \quad (4)$$

where
$$C_0 := 2C_{\pi^*}/(1-\gamma)$$
 and $C_1 := 2C_{\pi^*}/((1-\gamma)^3\mu_{tar, min})$.

The detailed proof of Proposition 2 is in Appendix B. The main insights are: (i) Similar to Proposition 1, the upper bound has two terms. Term (a) characterizes Q-NPG learning, while the suboptimality gap is mainly determined by the approximation error from $Q_{\rm src}$, equivalent to the crossdomain Bellman error (term (c)). (ii) Minimizing this error requires ϕ and ψ that reduce term (c). Motivated by Equation (4), we define cross-domain Bellman consistency.

Definition 4 (Cross-Domain Bellman Consistency). Let $\delta \geq 0$. A source-domain critic Q_{src} is said to be δ -Bellman-consistent under target domain policy π if there exist a pair of inter-domain mapping (ϕ, ψ) such that $\|\epsilon_{cd}(Q_{src}, \phi, \psi)\|_{d^{\pi}}$ is no more than δ .

Transferability of a Source-Domain Model. Given a source-domain critic $Q_{\rm src}$, if for any iteration t there exist inter-domain mappings $\phi^{(t)}$ and $\psi^{(t)}$ such that $Q_{\rm src}$ is δ -Bellman-consistent under $\pi^{(t)}$, then term (c) in (4) is bounded by $C_1\delta$. Thus, the transferability of a source-domain model is captured by δ . In the perfect transfer scenario, where source and target domains are identical and $Q_{\rm src}$ is optimal, setting ϕ and ψ as identity mappings ensures small δ for all t, yielding a small sub-optimality gap for sufficiently large T.

Limitations of DQT. By Proposition (2), a limitation of DQT is that with a poorly transferable source critic, the cross-domain Bellman error at each iteration t is large, so term (c) in (4) dominates the bound and prevents effective cross-domain transfer.

4.2 THE QAVATAR ALGORITHM

To address DQT's limitation, we propose QAvatar, which uses a hybrid critic consisting of a weighted combination of a learned target-domain Q function and a given source-domain Q function to enable reliable cross-domain knowledge transfer. This design allows QAvatar to improve sample efficiency in favorable scenarios while avoiding reliance on poorly transferable source models. Specifically, QAvatar comprises three major components:

• Inter-domain mapping: Under QAvatar, we propose to learn the inter-domain mappings $\phi: \mathcal{S}_{\text{tar}} \to \mathcal{S}_{\text{src}}$ and $\psi: \mathcal{A}_{\text{tar}} \to \mathcal{A}_{\text{src}}$ by minimizing the cross-domain Bellman loss as

$$\mathcal{L}_{\text{CD}}(\phi, \psi; Q_{\text{src}}, \pi_{\text{tar}}, \mathcal{D}_{\text{tar}}) := \hat{\mathbb{E}}_{(s, a, r_{\text{tar}}, s') \in \mathcal{D}_{\text{tar}}} \Big[\big| r_{\text{tar}} + \gamma \mathbb{E}_{a' \sim \pi_{\text{tar}}} [Q_{\text{src}}(\phi(s'), \psi(a'))] - Q_{\text{src}}(\phi(s), \psi(a)) \big| \Big],$$
(5)

Algorithm 2 QAvatar

Require: Source-domain Q function $Q_{\rm src}$.

- 1: Initialize the state mapping function ϕ , the action mapping function ψ , number of on-policy samples per iteration N_{tar} , the target-domain policy $\pi^{(0)}$, weight decay function $\alpha : \mathbb{N} \to [0, 1]$, and $\eta = (1 \gamma)\sqrt{1/T}$.
- 2: **for** iteration $t = 1, \dots, T$ **do**
- 3: Sample $\mathcal{D}_{tar}^{(t)} = \{(s, a, r, s')\}$ of $N_{tar}^{(t)}$ on-policy samples using $\pi^{(t)}$ in the target domain.
- 4: Update Q_{tar} by minimizing the TD loss in (1), i.e., $Q_{\text{tar}}^{(t)} \leftarrow \arg\min_{Q_{\text{tar}}} \mathcal{L}_{\text{TD}}(Q_{\text{tar}}; \pi^{(t)}, \mathcal{D}_{\text{tar}}^{(t)})$.
- 5: Update ϕ and ψ by minimizing (5), i.e., $\phi^{(t)}, \psi^{(t)} \leftarrow \arg\min_{\phi, \psi} \mathcal{L}_{CD}(\phi, \psi; Q_{src}, \pi^{(t)}, \mathcal{D}_{tar}^{(t)})$.
- 6: Defined weight parameter $\alpha(t) = \|\epsilon_{\mathrm{td}}^{(t)}\|_{\mathcal{D}_{\mathrm{ter}}^{(t)}} / (\|\epsilon_{\mathrm{cd}}(Q_{\mathrm{src}}, \phi^{(t)}, \psi^{(t)})\|_{\mathcal{D}_{\mathrm{ter}}^{(t)}} + \|\epsilon_{\mathrm{td}}^{(t)}\|_{\mathcal{D}_{\mathrm{ter}}^{(t)}}$
- 7: Update the target-domain policy by adapting NPG to CDRL as in (6).
- 8: end for
- 9: **Return** Target-domain policy $\pi_{\text{tar}}^{(T)} \sim \text{Uniform}(\{\pi^{(1)}, \cdots, \pi^{(T)}\})$.

where $Q_{\rm src}$ is the pre-trained source-domain Q function and $\mathcal{D}_{\rm tar} = \{(s,a,r_{\rm tar},s')\}$ denotes a set of target-domain samples drawn under $\pi_{\rm tar}$. Intuitively, the loss in (5) looks for a pair of mapping functions ϕ, ψ such that $Q_{\rm src}$ aligns as much with the target-domain transitions as possible.

- Target-domain Q function: To implement the hybrid critic, QAvatar maintains a target-domain Q function Q_{tar} , serving as the critic of the current target-domain policy. At each iteration t, Q_{tar} is obtained via policy evaluation by minimizing the TD loss $\mathcal{L}_{TD}(Q_{tar}; \pi_{tar}, \mathcal{D}_{tar})$, where $\mathcal{D}_{tar} = (s, a, r, s')$ are target-domain samples (Equation 1).
- NPG-like policy update with a weighted Q-function combination: QAvatar leverages both $Q_{\rm src}$ and $Q_{\rm tar}$ for policy updates. At each iteration t,

$$\pi^{(t+1)}(a|s) \propto \pi^{(t)}(a|s) \cdot \exp\left(\eta\left((1 - \alpha(t))Q_{\text{tar}}^{(t)}(s, a) + \alpha(t)Q_{\text{src}}(\phi^{(t)}(s), \psi^{(t)}(a))\right)\right), \quad (6)$$

where $\alpha:\mathbb{N}\to[0,1]$ is a weight function (see Section 4.3).

The pseudo code of QAvataris provided in Algorithm 2.

Remark 1. In line 6 of Algorithm 1 and line 8 of Algorithm 2, DQT and QAvatar output the final policy by selecting uniformly from all intermediate policies which is a standard procedure linking average sub-optimality to policy performance. In experiments, the last-iterate policy suffices and performs well.

4.3 THEORETICAL JUSTIFICATION OF QAVATAR

In this section, we present the theoretical result of QAvatar and thereby describe how to choose the proper decay parameter $\alpha(\cdot)$.

Definition 5 (Cross-Domain Action Value Function). For each state-action pair (s, a) and $t \in \mathbb{N}$, the cross-domain action value function $f^{(t)}(s, a)$ is defined as $f^{(t)}(s, a) := (1 - \alpha(t))Q_{tar}^{(t)}(s, a) + \alpha(t)Q_{src}(\phi^{(t)}(s), \psi^{(t)}(a))$.

We are ready to present the main theoretical result, and the detailed proof is provided in Appendix B.

Proposition 3. (Average Sub-Optimality) Under the QAvatar in Algorithm 2 and Assumption 1, the average sub-optimality over T iterations can be upper bounded as

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{s \sim \mu_{tar}} \left[V^{\pi^*}(s) - V^{\pi^{(t)}}(s) \right] \\
\leq \underbrace{\frac{\left[\log |\mathcal{A}_{tar}| + 1 \right]}{\sqrt{T}(1 - \gamma)}}_{(a)} + \underbrace{\frac{C_0}{T} \sum_{t=1}^{T} \mathbb{E}_{(s, a) \sim d^{\pi^{(t)}}} \left[\left| f^{(t)}(s, a) - Q^{\pi^{(t)}}(s, a) \right| \right]}_{(b)} \tag{7}$$

$$\leq \underbrace{\frac{\left[\log|\mathcal{A}_{tar}|+1\right]}{\sqrt{T}(1-\gamma)}}_{(a)} + \underbrace{\frac{C_{1}}{T}\sum_{t=1}^{T}\left(\alpha(t)\|\epsilon_{cd}(Q_{src},\phi^{(t)},\psi^{(t)})\|_{d^{\pi(t)}} + (1-\alpha(t))\|\epsilon_{td}^{(t)}\|_{d^{\pi(t)}}\right)}_{(c)}, \tag{8}$$

where $C_0 := 2C_{\pi^*}/(1-\gamma)$ and $C_1 := 2C_{\pi^*}/((1-\gamma)^3\mu_{tar, min})$.

Notably, the term (a) in (8) reflects the learning progress of NPG, and term (c) reflects the transferability of a source-domain critic $Q_{\rm src}$ and the error of policy evaluation for the target-domain policy.

A Hyperparameter-Free Design of $\alpha(t)$. Based on (8), for each iteration t, term (c) can be minimized by choosing $\alpha(t)$ as an indicator function, i.e., set to 1 when $\|\epsilon_{\rm cd}(Q_{\rm src},\phi^{(t)},\psi^{(t)})\|_{d^{\pi(t)}} < \|\epsilon_{\rm td}^{(t)}\|_{d^{\pi(t)}}$, and 0 otherwise. In practice, estimating the two error terms is noisy, so using an indicator can cause large fluctuations in $\alpha(t)$ and unstable training. To address this, we propose a smoother variant: $\alpha(t) = \|\epsilon_{\rm td}^{(t)}\|_{d^{\pi(t)}}/(\|\epsilon_{\rm cd}(Q_{\rm src},\phi^{(t)},\psi^{(t)})\|_{d^{\pi(t)}} + \|\epsilon_{\rm td}^{(t)}\|_{d^{\pi(t)}})$. Notably, this design is *hyperparameter-free* and incurs minimal deployment overhead.

Key Implications of Proposition 3: (1) Effective transfer lowers the upper bound of average suboptimality: In an ideal case with perfect mappings ϕ^*, ψ^* such that $L_{\text{CD}}(\phi^*, \psi^*; Q_{\text{src}}, \pi_{\text{tar}}, \mathcal{D}_{\text{tar}}) = 0$ for any π_{tar} , we obtain $\|\epsilon_{\text{cd}}(Q_{\text{src}}, \phi^*, \psi^*)\|_{d^{\pi_{\text{tar}}}} = 0$. Then $\alpha(t) = 1$ at all t, making term (c) in (8) vanish. The bound thus reduces to term (a), which becomes negligible as T grows. (2) QAvatar avoids being trapped by low-transfer critics. For a source critic only δ -Bellman-consistent with large δ , $\|\epsilon_{\text{cd}}(Q_{\text{src}}, \phi, \psi)\|_{d^{\pi^{(t)}}}$ remains large, so $\alpha(t) \approx 0$. Consequently, term (c) reduces to the standard TD error.

4.4 PRACTICAL IMPLEMENTATION OF QAVATAR

We extend the QAvatar framework in Algorithm 2 to a practical deep RL implementation. The pseudo code is provided in Algorithm 3 in Appendix.

- Learning the target-domain policy and the Q function. To go beyond the tabular setting, we extend QAvatar by connecting NPG with soft policy iteration (SPI) (Haarnoja et al., 2018). In the entropy-regularized RL setting, SPI is known to be a special case of NPG (Cen et al., 2022). Based on this connection, we choose to integrate QAvatar with soft actor-critic (SAC) (Haarnoja et al., 2018), *i.e.*, updating the target-domain critic Q_{tar} by the critic loss of SAC and updating the target-domain policy $\pi^{(t)}$ by the SAC policy loss with the weighted combination of Q_{tar} and Q_{src} of QAvatar.
- Learning the inter-domain mapping functions with an augmented flow model. Similar to the tabular setting, we learn inter-domain mappings by minimizing the cross-domain Bellman loss. In practical RL problems, state and action spaces are usually bounded, so the outputs of $\phi: \mathcal{S}_{tar} \to \mathcal{S}_{src}$ and $\psi: \mathcal{A}_{tar} \to \mathcal{A}_{src}$ must lie within feasible regions. As discussed in Section 2, adversarial learning is commonly used to address this (Taylor et al., 2008; Zhang et al., 2021; Gui et al., 2023; Zhu et al., 2024), but it can lead to unstable training. Therefore, we adopt the method of (Brahmanage et al., 2023), training a normalizing flow to map the outputs of the mapping functions into the feasible regions.

5 EXPERIMENTS

5.1 SETUP

Benchmark CDRL Methods. We compare QAvatar with recent CDRL benchmarks under different state-action spaces, including Cross-Morphology-Domain Policy Adaptation (CMD) (Gui et al., 2023), Cross-domain Adaptive Transfer (CAT) (You et al., 2022), and Policy Adaptation by Representation mismatch (PAR) (Lyu et al., 2024). For a fair comparison, all methods use the same source-domain models, including policy and corresponding Q-networks, pre-trained with SAC. We also evaluate both PPO-based CAT, the original version in (You et al., 2022), and SAC-based CAT. Notably, CMD is an enhanced version of (Zhang et al., 2021) that integrates dynamics cycle consistency to learn state-action correspondences.

To demonstrate sample efficiency, we also compare QAvatar with standard SAC (Haarnoja et al., 2018), which learns from scratch in the target domain, and with direct fine-tuning (FT) of the source models (Ha et al., 2024), equivalent to SAC with source feature initialization. Both serve as competitive baselines. Hyperparameters are provided in Appendix F.

Evaluation Environments.

- Locomotion: We use the standard MuJoCo environments, including Hopper-v3, HalfCheetah-v3 and Ant-v3, as the source domains and follow the same procedure as in (Zhang et al., 2021; Xu et al., 2023) to modify them for the target domains. The detailed morphologies are in Appendix F.
- Robot arm manipulation: We leverage Robosuite, a popular package for robot learning released by (Zhu et al., 2020) and evaluate our algorithm on door opening and table wiping. For each task, we use the Panda robot arm as the source domain and set the UR5e robot arm as the target domain.
- Goal Navigation: A natural transfer scenario occurs when the source and target domains share the same goal but differ in robot type. We use the Safety-Gym benchmark (Ray et al., 2019) and evaluate transfer from Car to Doggo, keeping the goal unchanged, specifically using CarGoal0 as the source and DoggoGoal0 as the target domain.

The dimensions of the state and action spaces of all the source-target pairs are in Table 3 in Appendix F. All the results reported below are averaged over 5 random seeds.

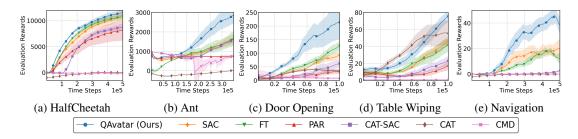


Figure 1: Training curves of QAvatar and benchmark methods: (a)-(b) Locomotion tasks; (d)-(e) Robot arm manipulation tasks in Robosuite; (f) Navigation task from CarGoal0 to DoggoGoal0.

5.2 EXPERIMENTAL RESULTS

Does QAvatar improve data efficiency?

Learning curves: As shown by Figure 1, we observe that QAvatar achieves improved data efficiency via cross-domain transfer than SAC throughout the training process in all the tasks, despite that these tasks have rather different dimensions as shown in Table 3.

CAT-SAC achieves moderate results on MuJoCo but transfers slowly to other tasks, as CAT-like methods lack guarantees and depend on parameter-based transfer, i.e., weighted combinations of source and target policy layers. Such methods assume shared feature representations (Zhuang et al., 2020), which often fails when domains differ. FT improves data efficiency over SAC on MuJoCo but learns slowly in Robosuite due to dissimilar state—action representations from different robot arms. CMD generally performs poorly and can be unstable (e.g., in Ant) owing to its adversarial mapping module. We attribute CMD's weakness to its unsupervised design, which ignores target-domain rewards.

Time to threshold: We provide Table 1 to mark the time to threshold. It shows that QAvatar requires only about 44% of the environment steps to achieve the threshold than SAC does in the best case.

Aggregated performance: To ensure a reliable comparison, we follow the guidelines of (Agarwal et al., 2021b) and calculate the interquartile mean (IQM) using rliable, which enables evaluation at an aggregated level. Figure 2 shows that QAvatarindeed achieves significantly better performance than all baselines.

Environment	Threshold	QAvatar	SAC	QAvatar / SAC
HalfCheetah	6000	126K	176K	0.71
Ant	1600	206K	346K	0.59
Door Opening	90	48K	98K	0.49
Table Wiping	45	72K	98K	0.73
Navigation	20	218K	490K	0.44

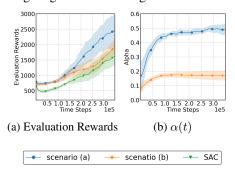
Table 1: Time to threshold of QAvatar and SAC

QAvatar CAT-SAC SAC CAT CMD PAR FT 0.25 0.50 0.75 Human Normalized Score

Figure 2: Aggregated IQMs (with 95% stratified bootstrap CIs) across tasks.

How does QAvatar perform under strong positive and negative transfer? We consider a task where the source domain is standard 'Ant-v3' and the target changes the goal to move backward, with all else unchanged. Here, $Q_{\rm src}$ and $Q_{\rm tar}$ are adversarial due to opposite goals. We evaluate QAvatar in two scenarios: (a) **Learning state/action mapping**: strong transferability exists, as Ant is symmetric along the front-back axis, allowing a perfect mapping. (b) **Fixing mapping as identity**: a strong negative transfer case, since $Q_{\rm src}$ provides adversarial reward signals. As shown in Figure 3, QAvatar captures both positive transfer (high $\alpha(t)$) and negative transfer (low $\alpha(t)$), demonstrating that $\alpha(t)$ reflects transferability.

Performance of QAvatar with a low-quality source domain: We evaluate this scenario in the Cheetah environment (Section 5.1) using a low-quality source model with a total return of 1000 (vs. \sim 7000 for the expert). Figure 4 illustrates the learning process and $\alpha(t)$ of QAvatar. Results show that when the source model is of low quality, $\alpha(t)$ decreases to a small value by the end of training, mitigating the effect of negative transfer.



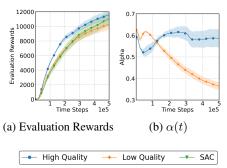


Figure 3: The training curve and the values of $\alpha(t)$ for QAvatar under strongly positive and strongly negative transfer scenarios.

Figure 4: The training curve and the values of $\alpha(t)$ in the Cheetah environment with a low-quality source model.

Does QAvatar still perform reliably well when the source and target with two unrelated transfer scenario? We evaluate transfer from original Hopper-v3 in MuJoCo to the table-wiping task in Robosuite. The configurations of these environments are provided in Section 5.1. Figure 6 shows that even when the source and target domains share no structural similarity, QAvatar still performs reliably and does not suffer from negative transfer.

How QAvatar perform on non-stationary environment? We use the Ant environment and introduce stochasticity by adding $\mathcal{N}(0,0.1)$ noise to rewards and $\mathcal{N}(0,0.05)$ to actions, following (Tessler et al., 2019). As shown in Figure 7, despite stochastic rewards and transitions, the inter-domain mapping is effectively learned, enabling positive transfer and faster learning in the target domain.

Extension: QAvatar with more than one source model. QAvatar can be readily extended for transfer from multiple source model. Similar to the idea of one source critic transfer, the weight $\alpha_i(t)$ for the *i*-th source critic $Q_{\mathrm{src},i}, \, \alpha_i(t) = (1/\|\epsilon_{\mathrm{cd}}(Q_{\mathrm{src},i},\phi_i^{(t)},\psi_i^{(t)})\|_{d^{\pi^{(t)}}})/(1/\|\epsilon_{\mathrm{td}}^{(t)}\|_{d^{\pi^{(t)}}} + \sum_{j=1}^N 1/\|\epsilon_{\mathrm{cd}}(Q_{\mathrm{src},j},\phi_i^{(t)},\psi_i^{(t)})\|_{d^{\pi^{(t)}}})$. Consider a two-source to one-target transfer scenario: (i) Source domain 1 (denoted by "src1") is Ant-v3 with the both front legs disabled; (ii) Source domain 2 (denoted by "src2") is Ant-v3 with the both back legs disabled. (iii) Target domain (denoted by "tar") is the original Ant-v3 with no modifications. Figure 8 shows QAvatarin multi-source cross-domain transfer can achieve higher transferability by leveraging the knowledge from two source domains.

6 CONCLUDING REMARKS

We propose cross-domain Bellman consistency as a measure of source-model transferability, and introduce QAvatar, the first CDRL method that reliably handles distinct state-action representations with performance guarantees. Using a hybrid critic and a hyperparameter-free weighting scheme, QAvatar achieves robust knowledge transfer even with weak source models. Experiments confirm its effectiveness for cross-domain RL. A limitation of our formulation is the assumption that target-domain data collection is costlier than training compute. Since QAvatar takes about twice the training time of SAC due to inter-domain mappings and the flow model, further acceleration would be needed when training efficiency is critical.

ETHICS STATEMENT

We conduct our research entirely in simulated environments, using no human participants or sensitive data. This work fully complies with the code of ethics.

REPRODUCIBILITY STATEMENT

The code for our experiments is provided in the supplementary material, along with a README file detailing the commands required to run the experiments. Furthermore, a comprehensive list of package dependencies is included to facilitate the recreation of the experimental environment.

USE OF LARGE LANGUAGE MODELS (LLMS)

Large language models (LLMs) were applied exclusively for linguistic refinement of the manuscript. No assistance was sought from LLMs in developing methods, performing experiments, or interpreting results.

BIBLIOGRAPHY

- Joshua Achiam. Spinning Up in Deep Reinforcement Learning. 2018.
- Alekh Agarwal, Nan Jiang, Sham M Kakade, and Wen Sun. Reinforcement learning: Theory and algorithms. *CS Dept.*, *UW Seattle, Seattle, WA, USA, Tech. Rep*, 32:96, 2019.
 - Alekh Agarwal, Mikael Henaff, Sham Kakade, and Wen Sun. Pc-pg: Policy cover directed exploration for provable policy gradient learning. In *Processing in Advances in Neural Information Processing Systems*, 2020.
 - Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learning Research*, 2021a.
 - Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C Courville, and Marc Bellemare. Deep reinforcement learning at the edge of the statistical precipice. In *Advances in Neural Information Processing Systems*, 2021b.
 - Haitham B Ammar, Karl Tuyls, Matthew E Taylor, Kurt Driessens, and Gerhard Weiss. Reinforcement learning transfer via sparse coding. In *International Conference on Autonomous Agents and Multiagent Systems*, 2012.
 - Haitham Bou Ammar and Matthew E Taylor. Reinforcement learning transfer via common subspaces. In *Adaptive and Learning Agents: International Workshop*, 2012.
 - Haitham Bou Ammar, Eric Eaton, Paul Ruvolo, and Matthew Taylor. Unsupervised cross-domain transfer in policy gradient reinforcement learning via manifold alignment. In *AAAI Conference on Artificial Intelligence*, 2015.
 - Paniz Behboudian, Yash Satsangi, Matthew E Taylor, Anna Harutyunyan, and Michael Bowling. Policy invariant explicit shaping: an efficient alternative to reward shaping. *Neural Computing and Applications*, 34(3):1673–1686, 2022.
 - Janaka Brahmanage, Jiajing Ling, and Akshat Kumar. Flowpg: Action-constrained policy gradient with normalizing flows. In *Advances in Neural Information Processing Systems*, 2023.
 - Shicong Cen, Chen Cheng, Yuxin Chen, Yuting Wei, and Yuejie Chi. Fast global convergence of natural policy gradient methods with entropy regularization. *Operations Research*, 2022.
 - Yevgen Chebotar, Ankur Handa, Viktor Makoviychuk, Miles Macklin, Jan Issac, Nathan Ratliff, and Dieter Fox. Closing the sim-to-real loop: Adapting simulation randomization with real world experience. In *IEEE International Conference on Robotics and Automation*, 2019.

- Dongsheng Ding, Kaiqing Zhang, Tamer Basar, and Mihailo Jovanovic. Natural policy gradient primal-dual method for constrained markov decision processes. In *Advances in Neural Information Processing Systems*, 2020.
 - Yuqing Du, Olivia Watkins, Trevor Darrell, Pieter Abbeel, and Deepak Pathak. Auto-tuned sim-to-real transfer. In *IEEE International Conference on Robotics and Automation*, 2021.
 - Mahidhar Dwarampudi and NV Reddy. Effects of padding on lstms and cnns. *arXiv preprint arXiv:1903.07288*, 2019.
 - Benjamin Eysenbach, Shreyas Chaudhari, Swapnil Asawa, Sergey Levine, and Ruslan Salakhutdinov. Off-dynamics reinforcement learning: Training for transfer with domain classifiers. In *International Conference on Learning Representations*, 2021.
 - Haiyuan Gui, Shanchen Pang, Shihang Yu, Sibo Qiao, Yufeng Qi, Xiao He, Min Wang, and Xue Zhai. Cross-domain policy adaptation with dynamics alignment. *Neural Networks*, 2023.
 - Abhishek Gupta, Coline Devin, YuXuan Liu, Pieter Abbeel, and Sergey Levine. Learning invariant feature spaces to transfer skills with reinforcement learning. In *International Conference on Learning Representations*, 2017.
 - Seokhyeon Ha, Sunbeom Jeong, and Jungwoo Lee. Domain-aware fine-tuning: Enhancing neural network adaptability. In *Association for the Advancement of Artificial Intelligence*, 2024.
 - Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning*, 2018.
 - Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In *International Conference on Machine Learning*, 2002.
 - Sham M Kakade. A natural policy gradient. In *Advances in Neural Information Processing Systems*, 2001.
 - Michail G. Lagoudakis and Ronald Parr. Model-free least-squares policy iteration. In *Advances in Neural Information Processing Systems*, 2001.
 - Alessandro Lazaric, Mohammad Ghavamzadeh, and Rémi Munos. Finite-sample analysis of least-squares policy iteration. *Journal of Machine Learning Research*, 2012.
 - Jinxin Liu, Zhang Hongyin, and Donglin Wang. DARA: Dynamics-Aware Reward Augmentation in Offline Reinforcement Learning. In *International Conference on Learning Representations*, 2022.
 - Jiafei Lyu, Chenjia Bai, Jing-Wen Yang, Zongqing Lu, and Xiu Li. Cross-domain policy adaptation by capturing representation mismatch. In *International Conference on Machine Learning*, 2024.
 - Steven Morad, Chris Lu, Ryan Kortvelesy, Stephan Liwicki, Jakob Foerster, and Amanda Prorok. Recurrent reinforcement learning with memoroids. In *Advances in Neural Information Processing Systems*, 2024.
 - Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 2009.
 - Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Sim-to-real transfer of robotic control with dynamics randomization. In *IEEE International Conference on Robotics and Automation*, 2018.
 - Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of Machine Learning Research*, 2021.
 - Aravind Rajeswaran, Sarvjeet Ghotra, Balaraman Ravindran, and Sergey Levine. Epopt: Learning robust neural network policies using model ensembles. In *International Conference on Learning Representations*, 2016.

- Alex Ray, Joshua Achiam, and Dario Amodei. Benchmarking safe exploration in deep reinforcement learning. *arXiv preprint arXiv:1910.01708*, 7(1):2, 2019.
 - Ram Ananth Sreenivasan, Hyun-Rok Lee, Yeonjeong Jeong, Jongseong Jang, Dongsub Shim, and Chi-Guhn Lee. A learnable similarity metric for transfer learning with dynamics mismatch. In *PRL Workshop Series Bridging the Gap Between AI Planning and Reinforcement Learning*, 2023.
 - Matthew E Taylor, Gregory Kuhlmann, and Peter Stone. Autonomous transfer for reinforcement learning. In *Autonomous Agents and Multiagent Systems*, 2008.
 - Chen Tessler, Yonathan Efroni, and Shie Mannor. Action robust reinforcement learning and applications in continuous control. In *International Conference on Machine Learning*, 2019.
 - Chang Wang and Sridhar Mahadevan. Manifold alignment without correspondence. In *International Joint Conference on Artificial Intelligence*, 2009.
 - Yue Wang, Yuting Liu, Wei Chen, Zhi-Ming Ma, and Tie-Yan Liu. Target transfer q-learning and its convergence analysis. *Neurocomputing*, 2020.
 - Karl Weiss, Taghi M Khoshgoftaar, and DingDing Wang. A survey of transfer learning. *Journal of Big Data*, 2016.
 - Kang Xu, Chenjia Bai, Xiaoteng Ma, Dong Wang, Bin Zhao, Zhen Wang, Xuelong Li, and Wei Li. Cross-domain policy adaptation via value-guided data filtering. In Advances in Neural Information Processing Systems, 2023.
 - Heng You, Tianpei Yang, Yan Zheng, Jianye Hao, and E. Taylor, Matthew. Cross-domain adaptive transfer reinforcement learning based on state-action correspondence. In *Uncertainty in Artificial Intelligence*, 2022.
 - Huizhen Yu and Dimitri P Bertsekas. Convergence results for some temporal difference methods based on least squares. *IEEE Transactions on Automatic Control*, 2009.
 - Rui Yuan, Simon S Du, Robert M Gower, Alessandro Lazaric, and Lin Xiao. Linear convergence of natural policy gradient methods with log-linear policies. In *International Conference on Learning Representations*, 2022.
 - Tom Zahavy, Matan Haroush, Nadav Merlis, Daniel J Mankowitz, and Shie Mannor. Learn what not to learn: Action elimination with deep reinforcement learning. In *Advances in Neural Information Processing Systems*, 2018.
 - Qiang Zhang, Tete Xiao, Alexei A Efros, Lerrel Pinto, and Xiaolong Wang. Learning cross-domain correspondence for control with dynamics cycle-consistency. In *International Conference on Learning Representations*, 2021.
 - Yifei Zhou, Ayush Sekhari, Yuda Song, and Wen Sun. Offline data enhanced on-policy policy gradient with provable guarantees. In *International Conference on Learning Representations*, 2024.
 - Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *International Conference on Computer Vision*, 2017.
 - Ruiqi Zhu, Tianhong Dai, and Oya Celiktutan. Cross domain policy transfer with effect cycle-consistency. In *IEEE International Conference on Robotics and Automation*, 2024.
 - Yuke Zhu, Josiah Wong, Ajay Mandlekar, Roberto Martín-Martín, Abhishek Joshi, Soroush Nasiriany, and Yifeng Zhu. Robosuite: A modular simulation framework and benchmark for robot learning. arXiv preprint arXiv:2009.12293, 2020.

APPENDICES

A SUPPORTING LEMMAS

Lemma 1 (Performance difference lemma). For any two policies π and π' , we have

$$V^{\pi'}(\mu) - V^{\pi}(\mu) = \frac{1}{1 - \gamma} \mathbb{E}_{s, a \sim d^{\pi'}} [A^{\pi}(s, a)],$$

 where $A^{\pi}(s,a) := Q^{\pi}(s,a) - V^{\pi}(s)$ is the advantage function.

Proof. This can be directly obtained from Lemma 6.1 in (Kakade & Langford, 2002). \Box

Lemma 2 ((Agarwal et al., 2019), Chapter 4). Let $\tau = (s_0, a_0, s_1, a_1, \cdots)$ denote the (random) trajectory generated under a policy π in an infinite-horizon MDP \mathcal{M} . For any function $f: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, we have

$$\mathbb{E}_{\tau} \left[\sum_{t=0}^{\infty} \gamma^t f(s_t, a_t) \right] = \frac{1}{1 - \gamma} \mathbb{E}_{(s, a) \sim d^{\pi}} \left[f(s, a) \right]. \tag{9}$$

Lemma 3 (Importance Ratio). Given a fixed policy π and a fixed state-action pair (s, a), let $p_k(s, a)$ denote the probability of reaching (s, a) under an initial distribution d^{π} and policy π after k time steps. Then, for any $k \in \mathbb{N}$, we have

$$\frac{p_k(s,a)}{d^{\pi}(s,a)} \le \frac{1}{(1-\gamma)\mu(s,a)}.$$
(10)

Proof. To begin with, recall the definition of d^{π} as

$$d^{\pi}(s,a) := (1 - \gamma) \Big(\mu(s,a) + \sum_{t=1}^{\infty} \gamma^t P(s_t = s, a_t = a; \pi, \mu) \Big) \equiv \sum_{t=0}^{\infty} \gamma^t P(s_t = s, a_t = a; \pi, \mu).$$

Let $s_{\text{next},k}$ and $a_{\text{next},k}$ denote the state and action after k time steps. Then, we can write down $p_k(s,a)$:

$$p_k(s, a) = \sum_{(s', a') \in S \times A} \mathbb{P}(s_{\text{next}, k} = s, a_{\text{next}, k} = a | s', a'; \pi) d^{\pi}(s', a')$$
(12)

$$= \sum_{(s',a')\in\mathcal{S}\times\mathcal{A}} \mathbb{P}(s_{\text{next},k} = s, a_{\text{next},k} = a|s',a';\pi) \cdot (1-\gamma) \cdot \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s', a_t = a';\pi,\mu)$$
(13)

$$= (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \sum_{s', a' \in \mathcal{S} \times \mathcal{A}} \mathbb{P}(s_{\text{next}, k} = s, a_{\text{next}, k} = a | s', a'; \pi, \mu) \cdot \mathbb{P}(s_t = s', a_t = a'; \pi, \mu)$$

$$= (1 - \gamma) \sum_{k=0}^{\infty} \gamma^{t} \mathbb{P}(s_{t+k} = s, a_{t+k} = a; \pi, \mu).$$
 (15)

Then, we have

$$\frac{p_k(s,a)}{d^{\pi}(s,a)} = \frac{(1-\gamma)\sum_{t=0}^{\infty} \gamma^t \, \mathbb{P}(s_{t+k} = s; a_{t+k} = a; \pi, \mu)}{(1-\gamma)\sum_{t=0}^{\infty} \gamma^t \, \mathbb{P}(s_t = s, a_t = a; \pi, \mu)}$$
(16)

$$= \frac{\sum_{t=0}^{\infty} \gamma^{t} \mathbb{P}(s_{t+k} = s, a_{t+k} = a; \pi, \mu)}{\sum_{t=0}^{\infty} \gamma^{t} \mathbb{P}(s_{t} = s, a_{t} = a; \pi, \mu)}$$
(17)

(14)

$$\leq \frac{\sum_{t=0}^{\infty} \gamma^{t}}{\sum_{t=0}^{\infty} \gamma^{t} \mathbb{P}(s_{t} = s; \pi, \mu)}$$

$$= \frac{1}{1} \cdot \frac{1}{\sum_{t=0}^{\infty} \sqrt{s}} (19)$$

where (18) holds by $\mathbb{P}(s_{t+k}=s,a_{t+k}=a;\pi,\mu)\leq 1$ and (19) holds by taking the sum of an infinite geometric sequence. By the fact that $\sum_{t=0}^{\infty}\gamma^t\,\mathbb{P}(s_t=s,a_t=a;\pi,\mu)=\mu(s,a)+\sum_{t=1}^{\infty}\gamma^t\,\mathbb{P}(s_t=s,a_t=a;\pi,\mu)$, we have

$$\frac{1}{1-\gamma} \cdot \frac{1}{\sum_{t=0}^{\infty} \gamma^{t} \, \mathbb{P}(s_{t}=s, a_{t}=a; \pi, \mu)} = \frac{1}{1-\gamma} \cdot \frac{1}{\mu(s, a) + \sum_{t=1}^{\infty} \gamma^{t} \, \mathbb{P}(s_{t}=s, a_{t}=a; \pi, \mu)}$$
(20)

$$\leq \frac{1}{(1-\gamma)\mu(s,a)} \tag{21}$$

where (21) holds by $\sum_{t=1}^{\infty} \gamma^t \mathbb{P}(s_t = s, a_t = a; \pi, \mu) \geq 0$.

Lemma 4. Let $\nu^{(t)}: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ and $\pi^{(t)}$ denote any tabular function used in the policy update and the policy at iteration t. That is,

$$\pi^{(t+1)}(a\mid s) \propto \pi^{(t)}(a\mid s) \exp\Big(\eta \nu^{(t)}(s,a)\Big).$$

Then, we assume that $\|\nu^{(t)}\|_{\infty} \leq 1/(1-\gamma)$ and setting learning rate $\eta = (1-\gamma)\sqrt{1/T}$ and optimal policy π^* , we have

$$\sum_{t=1}^{T} \mathbb{E}_{(s,a) \sim d^{\pi^*}} \left[\bar{\nu}^{(t)}(s,a) \right] \leq \frac{\sqrt{T} \left[\log |\mathcal{A}_{tar}| + 1 \right]}{1 - \gamma}$$

Proof. Let $\bar{\nu}^{(t)}(s,a) := \nu^{(t)}(s,a) - \nu^{(t)}(s,\pi^{(t)}(s))$. According to the policy update rule, at iteration t, the policy $\pi^{(t+1)}$ for the next iteration is updated by the formula:

$$\pi^{(t+1)}(a \mid s) = \frac{\pi^{(t)}(a \mid s) \exp\left(\eta \nu^{(t)}(s, a)\right)}{\sum_{a'} \pi^{(t)}\left(a' \mid s\right) \exp\left(\eta \nu^{(t)}(s, a')\right)} = \frac{\pi^{(t)}(a \mid s) \exp\left(\eta \bar{\nu}^{(t)}(s, a)\right)}{\sum_{a'} \pi^{(t)}\left(a' \mid s\right) \exp\left(\eta \bar{\nu}^{(t)}(s, a')\right)}. \tag{22}$$

Let $Z_t := \sum_{a'} \pi^{(t)} \left(a' \mid s \right) \exp \left(\eta \bar{\nu}^{(t)} \left(s, a' \right) \right)$. By multiplying both sides of (22) by Z_t , taking the logarithm, and then taking the expectation on both sides w.r.t $(s,a) \sim d^{\pi^*}$, we obtain

$$\mathbb{E}_{(s,a)\sim d^{\pi^*}} \left[\eta \bar{\nu}^{(t)}(s,a) \right] = \mathbb{E}_{(s,a)\sim d^{\pi^*}} \left[\log Z_t + \log \pi^{(t+1)}(a\mid s) - \log \pi^{(t)}(a\mid s) \right]. \tag{23}$$

Next, we bound the term $\log Z_t$. Note that $\eta \bar{\nu}^{(t)}(s,a) \leq \sqrt{1/T} \leq 1$ and the fact that $\exp(x) < 1 + x + x^2$ for any $x \leq 1$, we have

$$\log Z_t = \log \left(\sum_{a' \in A} \pi^{(t)} \left(a' \mid s \right) \exp \left(\eta \bar{\nu}^{(t)} \left(s, a' \right) \right) \right)$$

$$(24)$$

$$\leq \log \left(\sum_{a' \in \mathcal{A}} \pi^{(t)} \left(a' \mid s \right) \left[1 + \left(\eta \bar{\nu}^{(t)} \left(s, a' \right) \right) + \left(\eta \bar{\nu}^{(t)} \left(s, a' \right) \right)^{2} \right] \right) \tag{25}$$

$$\leq \log\left(1 + \frac{\eta^2}{(1-\gamma)^2}\right) \tag{26}$$

$$\leq \frac{\eta^2}{(1-\gamma)^2},\tag{27}$$

where (26) is because $\sum_{a'\in\mathcal{A}}\pi^{(t)}\left(a'\mid s\right)\bar{\nu}^{(t)}(s,a')=0$ and $\|\nu^{(t)}\|_{\infty}\leq 1/(1-\gamma)$, (27) is follow the fact that $\log(1+x)\leq x$ for any $x\geq 0$. Then, we have

$$\mathbb{E}_{(s,a)\sim d^{\pi^*}} \left[\eta \bar{\nu}^{(t)}(s,a) \right] \le \mathbb{E}_{(s,a)\sim d^{\pi^*}} \left[\log \pi^{(t+1)}(a\mid s) - \log \pi^{(t)}(a\mid s) + \frac{\eta^2}{(1-\gamma)^2} \right]. \tag{28}$$

By taking the summation over iterations on both sides of (28), we have

$$\sum_{t=1}^{T} \mathbb{E}_{(s,a) \sim d^{*}} \left[\eta \bar{\nu}^{(t)}(s,a) \right] \\
\leq \frac{T \eta^{2}}{(1-\gamma)^{2}} + \mathbb{E}_{(s,a) \sim d^{\pi^{*}}} \left[\log \pi^{(T+1)}(a \mid s) - \log \pi^{(1)}(a \mid s) \right].$$

Using the fact that $\log(\pi(a\mid s))\leq 0$ and $\pi^{(1)}(a\mid s)=\frac{1}{\mid\mathcal{A}_{\text{tar}}\mid}$, we have

$$\sum_{t=1}^T \mathbb{E}_{(s,a) \sim d^{\pi^*}} \left[\bar{\nu}^{(t)}(s,a) \right] \leq \frac{T\eta}{(1-\gamma)^2} + \frac{\log |\mathcal{A}_{\text{tar}}|}{\eta}.$$

By setting $\eta = (1 - \gamma)\sqrt{1/T}$, we have

$$\sum_{t=1}^T \mathbb{E}_{(s,a) \sim d^{\pi^*}} \left[\bar{\nu}^{(t)}(s,a) \right] \leq \frac{\sqrt{T} \left[\log |\mathcal{A}_{\text{tar}}| + 1 \right]}{1 - \gamma}$$

Lemma 5. Let $\nu^{(t)}: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ and $\pi^{(t)}$ denote value function used in the policy update and the policy at iteration t. That is,

$$\pi^{(t+1)}(a \mid s) \propto \pi^{(t)}(a \mid s) \exp\left(\eta \nu^{(t)}(s, a)\right).$$
 (29)

Then, we assume that $\|\nu^{(t)}\|_{\infty} \le 1/(1-\gamma)$ and setting learning rate $\eta = (1-\gamma)\sqrt{1/T}$ and optimal policy π^* , we have

$$\begin{split} &\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{s \sim \mu_{tar}} \Big[V^{\pi^*}(s) - V^{\pi^{(t)}}(s) \Big] \\ &\leq \frac{\left[\log |\mathcal{A}_{tar}| + 1 \right]}{\sqrt{T} (1 - \gamma)} + \frac{2C_{\pi^*}}{1 - \gamma} \frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{(s, a) \sim d^{\pi^{(t)}}} \left[\left| \nu^{(t)}(s, a) - Q^{\pi^{(t)}}(s, a) \right| \right] \end{split}$$

Proof.

$$V^{\pi^*}(\mu_{\text{tar}}) - V^{\pi^{(t)}}(\mu_{\text{tar}})$$

$$= \frac{1}{1 - \gamma} \mathbb{E}_{(s,a) \sim d_{\text{tar}}^{\pi^*}} \left[A^{\pi^{(t)}}(s,a) \right]$$
(30)

$$= \frac{1}{1 - \gamma} \mathbb{E}_{(s,a) \sim d_{\text{tar}}^{\pi^*}} \left[\bar{\nu}^{(t)}(s,a) - \bar{\nu}^{(t)}(s,a) + A^{\pi^{(t)}}(s,a) \right]$$
(31)

$$= \frac{1}{1 - \gamma} \mathbb{E}_{(s,a) \sim d_{\text{tar}}^{\pi^*}} \left[\bar{\nu}^{(t)}(s,a) \right] + \frac{1}{1 - \gamma} \mathbb{E}_{(s,a) \sim d_{\text{tar}}^{\pi^*}} \left[-\bar{\nu}^{(t)}(s,a) + A^{\pi^{(t)}}(s,a) \right]$$
(32)

$$\leq \frac{1}{1-\gamma} \mathbb{E}_{(s,a)\sim d_{\text{tar}}^{\pi^*}} \left[\bar{\nu}^{(t)}(s,a) \right] + \frac{1}{1-\gamma} \mathbb{E}_{(s,a)\sim d_{\text{tar}}^{\pi^*}} \left[\left| -\bar{\nu}^{(t)}(s,a) + A^{\pi^{(t)}}(s,a) \right| \right], \quad (33)$$

where (30) holds by the performance difference lemma (cf. Lemma 1), (31) is obtained by adding $^t(s,a) - \bar{\nu}^t(s,a)$, (32) is obtained by rearranging the terms in (31), and (33) holds by $x \leq |x|$, for all $x \in \mathbb{R}$. By the fact that $\|\frac{d^{\pi^*}}{d\pi^{(t)}}\|_{\infty} \leq C$, we have

$$\frac{1}{1-\gamma} \mathbb{E}_{(s,a)\sim d^{\pi^*}} \left[\bar{\nu}^{(t)}(s,a) \right] + \frac{1}{1-\gamma} \mathbb{E}_{s,a\sim d^{\pi^*}} \left[\left| -\bar{\nu}^{(t)}(s,a) + A^{\pi^{(t)}}(s,a) \right| \right] \\
\leq \frac{1}{1-\gamma} \mathbb{E}_{(s,a)\sim d^{\pi^*}} \left[\bar{\nu}^{(t)}(s,a) \right] + \frac{1}{1-\gamma} C \cdot \mathbb{E}_{s,a\sim d^{\pi^{(t)}}} \left[\left| -\bar{\nu}^{(t)}(s,a) + A^{\pi^{(t)}}(s,a) \right| \right] \tag{34}$$

(35)

Recall the definitions that $\bar{\nu}^{(t)}(s,a) := \nu^{(t)}(s,a) - \nu^{(t)}(s,\pi^{(t)}(s))$ and $A^{\pi^{(t)}}(s,a) := Q^{\pi^{(t)}}(s,a) - Q^{\pi^{(t)}}(s,\pi^{(t)}(s))$. Then, we have

$$\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\left|\bar{\nu}^{(t)}(s,a) - A^{\pi^{(t)}}(s,a)\right|\right] \\
= \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\left|\nu^{(t)}(s,a) - \nu^{(t)}(s,\pi^{(t)}(s)) - Q^{\pi^{(t)}}(s,a) + Q^{\pi^{(t)}}(s,\pi^{(t)}(s))\right|\right] \\
\leq \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\left|\nu^{(t)}(s,a) - Q^{\pi^{(t)}}(s,a)\right| + \left|Q^{\pi^{(t)}}(s,\pi^{(t)}(s)) - \nu^{(t)}(s,\pi^{(t)}(s))\right|\right] \tag{37}$$

where (37) holds by the fact that $|x+y| \le |x| + |y|$ for any $x, y \in \mathbb{R}$. Then, by linearity of expectation, we obtain

$$\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left| \nu^{(t)}(s,a) - Q^{\pi^{(t)}}(s,a) \right| + \left| Q^{\pi^{(t)}}(s,\pi^{(t)}(s)) - \nu^{(t)}(s,\pi^{(t)}(s)) \right| \right] \\
= \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left| \nu^{(t)}(s,a) - Q^{\pi^{(t)}}(s,a) \right| \right] + \mathbb{E}_{s\sim d^{\pi^{(t)}}} \left[\left| Q^{\pi^{(t)}}(s,\pi^{(t)}(s)) - \nu^{(t)}(s,\pi^{(t)}(s)) \right| \right] \\
= \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} 2 \left[\left| \nu^{(t)}(s,a) - Q^{\pi^{(t)}}(s,a) \right| \right] \tag{39}$$

where (39) holds by Jensen's inequality. Then, by substituting the result from (39) back into (34), we have

$$\frac{1}{1-\gamma} \mathbb{E}_{(s,a)\sim d^{\pi^*}} \left[\bar{\nu}^{(t)}(s,a) \right] + \frac{1}{1-\gamma} C \cdot \mathbb{E}_{s,a\sim d^{\pi(t)}} \left[\left| -\bar{\nu}^{(t)}(s,a) + A^{\pi^{(t)}}(s,a) \right| \right]$$

$$\leq \frac{1}{1-\gamma} \mathbb{E}_{(s,a)\sim d^{\pi^*}} \left[\bar{\nu}^{(t)}(s,a) \right] + \frac{2C}{1-\gamma} \cdot \mathbb{E}_{(s,a)\sim d^{\pi(t)}} \left[\left| \nu^{(t)}(s,a) - Q^{\pi^{(t)}}(s,a) \right| \right]$$
(41)

Next, summing over all iterations and combining with Lemma 4, we have

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{s \sim \mu_{tar}} \left[V^{\pi^*}(s) - V^{\pi^{(t)}}(s) \right]
\leq \frac{\left[\log |\mathcal{A}_{tar}| + 1 \right]}{\sqrt{T}(1 - \gamma)} + \frac{2C}{1 - \gamma} \frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{(s, a) \sim d^{\pi^{(t)}}} \left[\left| \nu^{(t)}(s, a) - Q^{\pi^{(t)}}(s, a) \right| \right]$$
(42)

Recall that for any policy π , we use d^{π} to denote the discounted state-action visitation distribution under policy π in the target domain.

Lemma 6. Under Algorithm 2, for any $t \in \mathbb{N}$, we have

$$\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\left|f^{t}(s,a) - Q^{\pi^{(t)}}(s,a)\right|\right] \\
\leq \frac{1}{(1-\gamma)^{2}\mu_{tar,min}}\left[(1-\alpha(t))\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\epsilon_{td}^{(t)}(s,a)\right] + \alpha(t)\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\epsilon_{cd}(s,a;Q_{src},\phi^{(t)},\psi^{(t)},\pi^{(t)})\right]\right] \\
\tag{43}$$

Proof. Recall the definition of $f^{(t)} := (1 - \alpha(t))Q_{\text{tar}}^{(t)}(s, a) + \alpha(t)Q_{\text{src}}(\phi^{(t)}(s), \psi^{(t)}(a))$, we have

$$\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\left|f^{(t)}(s,a) - Q^{\pi^{(t)}}(s,a)\right|\right]$$

$$= \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\left|(1-\alpha(t))Q_{\text{tar}}^{(t)}(s,a) + \alpha(t)Q_{\text{src}}(\phi^{(t)}(s),\psi^{(t)}(a)) - Q^{\pi^{(t)}}(s,a)\right|\right]$$

$$= \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\left|(1-\alpha(t))(Q_{\text{tar}}^{(t)}(s,a) - r_{\text{tar}}(s,a) + r_{\text{tar}}(s,a)) + r_{\text{tar}}(s,a)\right| + \alpha(t)(Q_{\text{src}}(\phi^{(t)}(s),\psi^{(t)}(a)) - r_{\text{tar}}(s,a) + r_{\text{tar}}(s,a)) - Q^{\pi^{(t)}}(s,a)\right]$$

$$(45)$$

$$\begin{aligned} & = \mathbb{E}_{(s,a) \sim d^{\pi(t)}} \left[\left| (1 - \alpha(t)) \left(Q_{\text{tar}}^{(t)}(s,a) - r_{\text{tar}}(s,a) + r_{\text{tar}}(s,a) - \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[Q_{\text{tar}}^{(t)}(s',a')] \right. \\ & + \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[Q_{\text{tar}}^{(t)}(s',a')] \right) + \alpha(t) \left(Q_{\text{src}}(\phi^{(t)}(s), \psi^{(t)}(a)) - r_{\text{tar}}(s,a) + r_{\text{tar}}(s,a) \right. \\ & + \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[Q_{\text{tar}}(s',a')] \right) + \alpha(t) \left(Q_{\text{src}}(\phi^{(t)}(s), \psi^{(t)}(a)) - r_{\text{tar}}(s,a) + r_{\text{tar}}(s,a) \right. \\ & - \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[Q_{\text{src}}(\phi^{(t)}(s'), \psi^{(t)}(a'))] + \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[Q_{\text{src}}(\phi^{(t)}(s'), \psi^{(t)}(a'))] \right) \\ & - \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[Q_{\text{src}}(\phi^{(t)}(s'), \psi^{(t)}(a'))] + \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[Q_{\text{src}}(\phi^{(t)}(s'), \psi^{(t)}(a'))] \right) \\ & - \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[Q_{\text{src}}(\phi^{(t)}(s'), \psi^{(t)}(a'))] + \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[Q_{\text{tar}}(s',a')] \right) \\ & - \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[Q_{\text{src}}(\phi^{(t)}(s'), \psi^{(t)}(a'))] \right) \\ & - \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[Q_{\text{tar}}(s',a') - r_{\text{tar}}(s,a) - \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[Q_{\text{src}}(\phi^{(t)}(s'), \psi^{(t)}(a'))] \right) \\ & - \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[Q_{\text{src}}(\phi^{(t)}(s'), \psi^{(t)}(a'))] \right) \\ & + \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[f^{(t)}(s',a')] + r_{\text{tar}}(s,a) - \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[Q_{\text{tar}}(s',a')] \right) \\ & - \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[f^{(t)}(s',a')] + r_{\text{tar}}(s,a) - \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[Q_{\text{tar}}(s',a')] \right) \\ & - \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[f^{(t)}(s',a')] + r_{\text{tar}}(s,a) - \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[Q_{\text{tar}}(s',a')] \right) \\ & - \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[f^{(t)}(s',a')] + r_{\text{tar}}(s,a) - Q^{\pi^{(t)}}(s,a) \right] \right], \\ & + \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[f^{(t)}(s',a')] + r_{\text{tar}}(s,a) - Q^{\pi^{(t)}}(s,a) \right] \right], \\ & + \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[f^{(t)}(s',a')] + r_{\text{tar}}(s,a) - Q^{\pi^{(t)}}(s,a) \right] \right], \\ & + \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[f^{(t)}(s',a')] + r_{\text{tar}}(s,a) - Q^{\pi^{(t)}}(s,a) \right] \right], \\ & + \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[f^{(t)}(s',a')] + r_{\text{t$$

where we obtain (45) by adding the dummy terms $(1-\alpha(t))(-r_{\text{tar}}(s,a)+r_{\text{tar}}(s,a))$ and $\alpha(t)(-r_{\text{tar}}(s,a)+r_{\text{tar}}(s,a))$ to the inner part of (44), (46) is obtained by adding $(1-\alpha(t))(-\gamma\mathbb{E}_{s'\sim P_{\text{tar}}(\cdot|s,a)}[Q_{\text{tar}}^{(t)}(s',a')]+\gamma\mathbb{E}_{s'\sim P_{\text{tar}}(\cdot|s,a)}[Q_{\text{tar}}^{(t)}(s',a')])$ and $\alpha(t)(-\alpha(t))$

$$\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left| (1-\alpha(t)) \left(Q_{\text{tar}}^{(t)}(s,a) - r_{\text{tar}}(s,a) - \gamma \mathbb{E}_{s'\sim P_{\text{tar}}(\cdot|s,a)} [Q_{\text{tar}}^{(t)}(s',a')] \right) \right. \\
\left. + \alpha(t) \left(Q_{\text{src}}(\phi^{(t)}(s), \psi^{(t)}(a)) - r_{\text{tar}}(s,a) - \gamma \mathbb{E}_{s'\sim P_{\text{tar}}(\cdot|s,a)} [Q_{\text{src}}(\phi^{(t)}(s'), \psi^{(t)}(a'))] \right) \\
\left. a' \sim \pi^{(t)}(\cdot|s') \right. \\
\left. + \gamma \mathbb{E}_{s'\sim P_{\text{tar}}(\cdot|s,a)} [f^{(t)}(s',a')] + r_{\text{tar}}(s,a) - Q^{\pi^{(t)}}(s,a) \\
\left. a' \sim \pi^{(t)}(\cdot|s') \right. \\
\left. + \gamma \mathbb{E}_{s''\sim P_{\text{tar}}(\cdot|s,a)} [Q^{\pi^{(t)}}(s'',a'')] - \gamma \mathbb{E}_{s''\sim P_{\text{tar}}(\cdot|s,a)} [Q^{\pi^{(t)}}(s'',a'')] \right| \right] \\
\left. a'' \sim \pi^{(t)}(\cdot|s'') \right. \\$$

$$\leq \mathbb{E}_{(s,a) \sim d^{\pi(t)}} \left[\left| (1 - \alpha(t)) \left(Q_{\text{tur}}^{(t)}(s, a) - r_{\text{tar}}(s, a) - \gamma \mathbb{E}_{s' \sim P_{\text{tur}}(\cdot \mid s, a)} \left[Q_{\text{tar}}^{(t)}(s', a') \right] \right] \right| \\
+ \left| \alpha(t) \left(Q_{\text{src}}(\phi^{(t)}(s), \psi^{(t)}(a)) - r_{\text{tar}}(s, a) - \gamma \mathbb{E}_{s' \sim P_{\text{tur}}(\cdot \mid s, a)} \left[Q_{\text{src}}(\phi^{(t)}(s'), \psi^{(t)}(a')) \right] \right) \right| \\
+ \left| \gamma \mathbb{E}_{s' \sim P_{\text{tur}}(\cdot \mid s, a)} \left[f^{(t)}(s', a') \right] + r_{\text{tar}}(s, a) - Q^{\pi^{(t)}}(s, a) \right| \\
+ \left| \gamma \mathbb{E}_{s'' \sim P_{\text{tur}}(\cdot \mid s, a)} \left[Q^{\pi^{(t)}}(s'', a'') \right] - \gamma \mathbb{E}_{s'' \sim P_{\text{tur}}(\cdot \mid s, a)} \left[Q^{\pi^{(t)}}(s'', a'') \right] \right] \right] \\
\leq \mathbb{E}_{(s,a) \sim d^{\pi^{(t)}}} \left[\left(1 - \alpha(t) \right) \left[Q_{\text{tar}}^{(t)}(s, a) - r_{\text{tar}}(s, a) - \gamma \mathbb{E}_{s' \sim P_{\text{tur}}(\cdot \mid s, a)} \left[Q_{\text{tar}}^{(t)}(s', a') \right] \right] \\
= \varepsilon_{\text{tot}}^{(t)}(s, a) \\
+ \alpha(t) \left[\left(Q_{\text{src}}(\phi^{(t)}(s), \psi^{(t)}(a)) - r_{\text{tar}}(s, a) - \gamma \mathbb{E}_{s' \sim P_{\text{tur}}(\cdot \mid s, a)} \left[Q_{\text{src}}(\phi^{(t)}(s'), \psi^{(t)}(a')) \right] \right) \right| \\
= \varepsilon_{\text{tot}}^{(t)}(s, a) \\
+ \alpha(t) \left[\left(Q_{\text{src}}(\phi^{(t)}(s), \psi^{(t)}(a)) - r_{\text{tar}}(s, a) - \gamma \mathbb{E}_{s' \sim P_{\text{tur}}(\cdot \mid s, a)} \left[Q_{\text{src}}(\phi^{(t)}(s'), \psi^{(t)}(a')) \right] \right) \right| \\
= \varepsilon_{\text{tot}}^{(t)}(s, a) \\
+ \left| \gamma \mathbb{E}_{s' \sim P_{\text{tur}}(\cdot \mid s, a)} \left[f^{(t)}(s', a') \right] - \gamma \mathbb{E}_{s'' \sim P_{\text{tur}}(\cdot \mid s, a)} \left[Q^{\pi^{(t)}}(s'', a'') \right] \right| \\
= \varepsilon_{\text{tot}}^{(t)}(s, a') \\
=$$

where (50) holds by triangle inequality, (51) holds by the facts that $0 \leq \alpha(t) \leq 1$ and $0 \leq 1 - \alpha(t) \leq 1$, (52) holds by coupling (s',a') and (s'',a'') and applying Bellman expectation equation as well as the definitions that $\epsilon_{\rm td}^{(t)}(s,a) := |Q_{\rm tar}^{(t)}(s,a) - r_{\rm tar}(s,a) - \gamma \mathbb{E}_{s' \sim P_{\rm tar}(\cdot|s,a)}[Q_{\rm tar}^{(t)}(s',a')]|$ and $\epsilon_{\rm cd}(s,a;Q_{\rm src},\phi^{(t)},\psi^{(t)},\pi^{(t)}) := a' \sim \pi^{(t)}(\cdot|s')$ $|Q_{\rm src}(\phi^{(t)}(s),\psi^{(t)}(a))|$. By recursively aparallow $|Q_{\rm src}(\phi^{(t)}(s),\psi^{(t)}(a))|$.

plying the procedure from (44) to (52) to $\left|f^{(t)}(s',a') - Q^{\pi^{(t)}}(s',a')\right|$, we obtain a bound on $\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\left(f^{(t)}(s,a) - Q^{\pi^{(t)}}(s,a)\right)^2\right]$ as follows:

$$\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\left|f^{(t)}(s,a) - Q^{\pi^{(t)}}(s,a)\right|\right] \\
\leq \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\left|\left(1 - \alpha(t)\right)\epsilon_{\mathsf{td}}^{(t)}(s,a) + \alpha(t)\epsilon_{\mathsf{cd}}(s,a;Q_{\mathsf{src}},\phi^{(t)},\psi^{(t)},\pi^{(t)}) \\
+ \gamma \mathbb{E}_{s'\sim P_{\mathsf{tar}}(\cdot|s,a)}\left[\left|f^{(t)}(s',a') - Q^{\pi^{(t)}}(s',a')\right|\right]\right|\right]$$
(53)

$$\leq \mathbb{E}_{(s,a)\sim d^{\pi(t)}} \left[\left| (1-\alpha(t))\epsilon_{td}^{(t)}(s,a) + \alpha(t)\epsilon_{cd}(s,a;Q_{src},\phi^{(t)},\psi^{(t)},\pi^{(t)}) \right. \\
+ \gamma \mathbb{E}_{s'\sim P_{tar}(\cdot|s,a)} \left[(1-\alpha(t))\epsilon_{td}^{(t)}(s',a') + \alpha(t)\epsilon_{cd}(s',a';Q_{src},\phi^{(t)},\psi^{(t)},\pi^{(t)}) \right. \\
+ \gamma \mathbb{E}_{s''\sim P_{tar}(\cdot|s',a')} \left[\left| f^{(t)}(s'',a'') - Q^{\pi^{(t)}}(s'',a'') \right| \right] \right] \right] \\
\leq \mathbb{E}_{(s,a)\sim d^{\pi(t)}} \left[\left| (1-\alpha(t))\epsilon_{td}^{(t)}(s,a) + \alpha(t)\epsilon_{cd}(s,a;Q_{src},\phi^{(t)},\psi^{(t)},\pi^{(t)}) \right. \\
+ \frac{1}{(1-\gamma)\mu_{tar,min}} \left(\gamma(1-\alpha(t))\epsilon_{td}^{(t)}(s,a) + \gamma\alpha(t)\epsilon_{cd}(s,a;Q_{src},\phi^{(t)},\psi^{(t)},\pi^{(t)}) \right. \\
+ \gamma^{2} (1-\alpha(t))\epsilon_{td}^{(t)}(s,a) + \gamma^{2}\alpha(t)\epsilon_{cd}(s,a;Q_{src},\phi^{(t)},\psi^{(t)},\pi^{(t)}) + \cdots \right) \right] \\
\leq \frac{1}{(1-\gamma)^{2}\mu_{tar,min}} \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left| (1-\alpha(t))\epsilon_{td}^{(t)}(s,a) + \alpha(t)\epsilon_{cd}(s,a;Q_{src},\phi^{(t)},\psi^{(t)},\pi^{(t)}) \right| \right] (56) \\
= \frac{1}{(1-\gamma)^{2}\mu_{tar,min}} \left[(1-\alpha(t))\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\epsilon_{td}^{(t)}(s,a) \right] + \alpha(t)\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\epsilon_{cd}(s,a;Q_{src},\phi^{(t)},\psi^{(t)},\pi^{(t)}) \right] \right] (57)$$

where (54) holds by applying the procedure from (44) to (52) to $f^{(t)}(s',a') - Q^{\pi^{(t)}}(s',a')$, (55) holds by applying the procedure from (44) to (52) to all the subsequent time steps and using importance sampling with the importance ratio bound in Lemma 3 and then using the same dummy variables (s,a) for all the subsequent state-action pairs, (57) holds by taking the sum of an infinite geometric sequence.

B PROOFS OF THE PROPOSITIONS

We first present the proof of Proposition 3 in Appendix B.1 and then establish Proposition 2 and 1 by a similar argument in Appendix B.3.

B.1 PROOF OF PROPOSITION 3

Proposition 3. (Average Sub-Optimality) Under the QAvatar in Algorithm 2 and Assumption 1, the average sub-optimality over T iterations can be upper bounded as

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{s \sim \mu_{tar}} \left[V^{\pi^*}(s) - V^{\pi^{(t)}}(s) \right] \\
\leq \underbrace{\frac{\left[\log |\mathcal{A}_{tar}| + 1 \right]}{\sqrt{T}(1 - \gamma)}}_{(a)} + \underbrace{\frac{C_0}{T} \sum_{t=1}^{T} \mathbb{E}_{(s, a) \sim d^{\pi^{(t)}}} \left[\left| f^{(t)}(s, a) - Q^{\pi^{(t)}}(s, a) \right| \right]}_{(b)} \\
\leq \underbrace{\frac{\left[\log |\mathcal{A}_{tar}| + 1 \right]}{\sqrt{T}(1 - \gamma)}}_{(a)} + \underbrace{\frac{C_1}{T} \sum_{t=1}^{T} \left(\alpha(t) \|\epsilon_{cd}(Q_{src}, \phi^{(t)}, \psi^{(t)})\|_{d^{\pi^{(t)}}} + (1 - \alpha(t)) \|\epsilon_{td}^{(t)}\|_{d^{\pi^{(t)}}} \right)}_{(c)}, \tag{8}$$

where $C_0 := 2C_{\pi^*}/(1-\gamma)$ and $C_1 := 2C_{\pi^*}/((1-\gamma)^3\mu_{tar, min})$.

Proof. Using Lemma 5 and setting $v^{(t)} = f^{(t)}$, we have

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{s \sim \mu_{tar}} \left[V^{\pi^*}(s) - V^{\pi^{(t)}}(s) \right] \\
\leq \frac{\left[\log |\mathcal{A}_{tar}| + 1 \right]}{\sqrt{T}(1 - \gamma)} + \frac{2C}{1 - \gamma} \frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{(s,a) \sim d^{\pi^{(t)}}} \left[\left| f^{(t)}(s,a) - Q^{\pi^{(t)}}(s,a) \right| \right] \tag{58}$$

This establishes the first inequality. Furthermore, recall the definitions of $\epsilon_{\rm td}^{(t)}(s,a)$ and $\epsilon_{\rm cd}(s,a;Q_{\rm src},\phi,\psi,\pi)$ as

$$\epsilon_{\rm td}^{(t)}(s,a) := |Q_{\rm tar}^{(t)}(s,a) - r_{\rm tar}(s,a) - \gamma \mathbb{E}_{s' \sim P_{\rm tar}(\cdot|s,a)}[Q_{\rm tar}^{(t)}(s',a')]|,$$

$$\epsilon_{\rm td}^{(t)}(s,a) := |Q_{\rm tar}^{(t)}(s,a) - r_{\rm tar}(s,a) - \gamma \mathbb{E}_{s' \sim P_{\rm tar}(\cdot|s,a)}[Q_{\rm tar}^{(t)}(s',a')]|,$$
(59)

$$\epsilon_{\mathrm{cd}}(s, a; Q_{\mathrm{src}}, \phi, \psi, \pi) := \left| Q_{\mathrm{src}}(\phi(s), \psi(a)) - r_{\mathrm{tar}}(s, a) - \gamma \mathbb{E}_{s' \sim P_{\mathrm{tar}}(\cdot \mid s, a), a' \sim \pi(\cdot \mid s')} [Q_{\mathrm{src}}(\phi(s'), \psi(a'))] \right|. \tag{60}$$

We also define the weighted ℓ_1 norm under state-action distribution induced by any policy π as

$$\|\epsilon_{\mathsf{td}}^{(t)}\|_{d^{\pi}} := \mathbb{E}_{(s,a) \sim d^{\pi}} \left[\epsilon_{\mathsf{td}}^{(t)}(s,a)\right],\tag{61}$$

(64)

$$\|\epsilon_{\rm cd}(Q_{\rm src}, \phi^{(t)}, \psi^{(t)})\|_{d^{\pi}} := \mathbb{E}_{(s,a) \sim d^{\pi}} \left[\epsilon_{\rm cd}(s, a; Q_{\rm src}, \phi^{(t)}, \psi^{(t)}, \pi)\right]. \tag{62}$$

For the second inequality, by Lemma 6, we have

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{s \sim \mu_{\text{tar}}} \left[V^{\pi^*}(s) - V^{\pi^{(t)}}(s) \right] \\
\leq \frac{\left[\log |\mathcal{A}_{\text{tar}}| + 1 \right]}{\sqrt{T} (1 - \gamma)} + \frac{2C}{1 - \gamma} \frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{(s, a) \sim d^{\pi^{(t)}}} \left[\left| f^{(t)}(s, a) - Q^{\pi^{(t)}}(s, a) \right| \right] \\
\leq \frac{\left[\log |\mathcal{A}_{\text{tar}}| + 1 \right]}{\sqrt{T} (1 - \gamma)} + \frac{2C}{(1 - \gamma)^3 \mu_{\text{tar,min}}} \frac{1}{T} \sum_{t=1}^{T} \left[(1 - \alpha(t)) \|\epsilon_{\text{td}}^{(t)}\|_{d^{\pi^{(t)}}} + \alpha(t) \|\epsilon_{\text{cd}}(Q_{\text{src}}, \phi^{(t)}, \psi^{(t)})\|_{d^{\pi^{(t)}}} \right]$$
(63)

This completes the proof of Proposition 3. Additionally, by choosing $\alpha(t) = \frac{\|\epsilon_{\rm id}^{(t)}\|_{d^{\pi(t)}}}{\|\epsilon_{\rm cd}(Q_{\rm src},\phi^{(t)},\psi^{(t)})\|_{d^{\pi(t)}} + \|\epsilon_{\rm id}^{(t)}\|_{d^{\pi(t)}}}$ (as discussed in Section 4), we have

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{s \sim \mu_{\text{tar}}} \left[V^{\pi^*}(s) - V^{\pi^{(t)}}(s) \right] \\
\leq \frac{2}{(1-\gamma)^2} \sqrt{\frac{\log(\mathcal{A}_{\text{tar}})}{T}} + \frac{4\sqrt{2C}}{(1-\gamma)^3 \mu_{\text{tar,min}}} \frac{1}{T} \sum_{t=1}^{T} \frac{\|\epsilon_{\text{cd}}(Q_{\text{src}}, \phi^{(t)}, \psi^{(t)})\|_{d^{\pi^{(t)}}} \cdot \|\epsilon_{\text{td}}^{(t)}\|_{d^{\pi^{(t)}}}}{\|\epsilon_{\text{cd}}(Q_{\text{src}}, \phi^{(t)}, \psi^{(t)})\|_{d^{\pi^{(t)}}} + \|\epsilon_{\text{td}}^{(t)}\|_{d^{\pi^{(t)}}}}.$$
(65)

B.2 Proof of Proposition 2

Proposition 2. Under the DQT method in Algorithm 1 and Assumption 1, the average sub-optimality over T iterations is upper bounded as

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{s \sim \mu_{tar}} \left[V^{\pi^*}(s) - V^{\pi^{(t)}}(s) \right] \leq \underbrace{\frac{\left[\log |\mathcal{A}_{tar}| + 1 \right]}{\sqrt{T}(1 - \gamma)}}_{(a)} + \underbrace{\frac{C_0}{T} \sum_{t=1}^{T} \left\| \left| Q_{src}(\phi^{(t)}, \psi^{(t)}) - Q^{\pi^{(t)}} \right| \right\|_{d^{\pi^{(t)}}}}_{(b)} \\
\leq \underbrace{\frac{\left[\log |\mathcal{A}_{tar}| + 1 \right]}{\sqrt{T}(1 - \gamma)}}_{(a)} + \underbrace{\frac{C_1}{T} \sum_{t=1}^{T} \left\| \epsilon_{cd}(Q_{src}, \phi^{(t)}, \psi^{(t)}) \right\|_{d^{\pi^{(t)}}}}_{(c)}, \quad (4)$$

where $C_0 := 2C_{\pi^*}/(1-\gamma)$ and $C_1 := 2C_{\pi^*}/((1-\gamma)^3\mu_{tar, min})$.

Proof. Notably, since the Proposition 2 is a special case of Proposition 3, we can simply follow all the steps taken for Proposition 3 and set $\alpha(t)=1$ for all t to establish Proposition 2. More specifically, we can replace $f^{(t)}(s,a)$ with $Q_{\rm src}(\phi^{(t)}(s),\psi^{(t)}(a))$. Accordingly, under $\alpha(t)=1$ for all t, Lemma 6 can be simply rewritten as

$$\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\left|Q_{\text{src}}(\phi^{(t)}(s),\psi^{(t)}(a)) - Q^{\pi^{(t)}}(s,a)\right|\right]$$
(66)

$$\leq \frac{1}{(1-\gamma)^{2}\mu_{\text{tar,min}}} \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\epsilon_{\text{cd}}(s,a;Q_{\text{src}},\phi^{(t)},\psi^{(t)},\pi^{(t)}) \right]. \tag{67}$$

Similarly, Lemma 5 can be be rewritten as

$$\begin{split} &\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{s \sim \mu_{\text{tar}}} \Big[V^{\pi^*}(s) - V^{\pi^{(t)}}(s) \Big] \\ &\leq \frac{\left[\log |\mathcal{A}_{\text{tar}}| + 1 \right]}{\sqrt{T} (1 - \gamma)} + \frac{2C_{\pi^*}}{1 - \gamma} \frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{(s, a) \sim d^{\pi^{(t)}}} \left[\left| Q_{\text{src}}(\phi^{(t)}(s), \psi^{(t)}(a)) - Q^{\pi^{(t)}}(s, a) \right| \right] \end{split}$$

From the combination of the two results.

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{s \sim \mu_{\text{tar}}} \Big[V^{\pi^*}(s) - V^{\pi^{(t)}}(s) \Big] \leq \frac{\left[\log |\mathcal{A}_{\text{tar}}| + 1 \right]}{\sqrt{T} (1 - \gamma)} + \frac{2C_{\pi^*}}{(1 - \gamma)^3 \mu_{\text{tar,min}} T} \sum_{t=1}^{T} \left\| \epsilon_{\text{cd}}(Q_{\text{src}}, \phi^{(t)}, \psi^{(t)}) \right\|_{d^{\pi^{(t)}}}. \tag{68}$$

B.3 Proof of Proposition 1

Proposition 1. Under the tabular and approximate-Q settings, and Assumption 1, the average sub-optimality of Q-NPG over T iterations is upper bounded by

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{s \sim \mu_{tar}} \left[V^{\pi^*}(s) - V^{\pi^{(t)}}(s) \right] \\
\leq \underbrace{\frac{\left[\log |\mathcal{A}_{tar}| + 1 \right]}{\sqrt{T}(1 - \gamma)}}_{(a)} + \underbrace{\frac{C_0}{T} \sum_{t=1}^{T} \left\| \left| Q_{tar}^{(t)} - Q^{\pi^{(t)}} \right| \right\|_{d^{\pi^{(t)}}}}_{(b)} \leq \underbrace{\frac{\left[\log |\mathcal{A}_{tar}| + 1 \right]}{\sqrt{T}(1 - \gamma)}}_{(a)} + \underbrace{\frac{C_1}{T} \sum_{t=1}^{T} \left\| \epsilon_{td}^{(t)} \right\|_{d^{\pi^{(t)}}}}_{(c)}, \quad (2)$$

where $C_0 := 2C_{\pi^*}/(1-\gamma)$ and $C_1 := 2C_{\pi^*}/((1-\gamma)^3\mu_{tar, min})$.

Proof. Notably, since the Proposition 1 is a special case of Proposition 3, we can simply follow all the steps taken for Proposition 3 and set $\alpha(t)=0$ for all t to establish Proposition 1. More specifically, we can replace $f^{(t)}(s,a)$ with $Q^{(t)}_{tar}(s,a)$. Accordingly, under $\alpha(t)=0$ for all t, Lemma 6 can be simply rewritten as

$$\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\left|Q_{\text{tar}}^{(t)}(s,a) - Q^{\pi^{(t)}}(s,a)\right|\right]$$
 (69)

$$\leq \frac{1}{(1-\gamma)^2 \mu_{\text{tar,min}}} \mathbb{E}_{(s,a) \sim d^{\pi^{(t)}}} \left[\epsilon_{\text{td}}(s,a) \right]. \tag{70}$$

Similarly, Lemma 5 can be be rewritten as

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{s \sim \mu_{\text{tar}}} \left[V^{\pi^*}(s) - V^{\pi^{(t)}}(s) \right] \\
\leq \frac{\left[\log |\mathcal{A}_{\text{tar}}| + 1 \right]}{\sqrt{T}(1 - \gamma)} + \frac{2C_{\pi^*}}{1 - \gamma} \frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{(s, a) \sim d^{\pi^{(t)}}} \left[\left| Q_{\text{tar}}^{(t)}(s, a) - Q^{\pi^{(t)}}(s, a) \right| \right]$$

From the combination of the two results,

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{s \sim \mu_{\text{tar}}} \left[V^{\pi^*}(s) - V^{\pi^{(t)}}(s) \right] \leq \frac{\left[\log |\mathcal{A}_{\text{tar}}| + 1 \right]}{\sqrt{T} (1 - \gamma)} + \frac{2C_{\pi^*}}{(1 - \gamma)^3 \mu_{\text{tar,min}} T} \sum_{t=1}^{T} \left\| \epsilon_{\text{td}} \right\|_{d^{\pi^{(t)}}}. \tag{71}$$

C A DETAILED DESCRIPTION OF RELATED WORK

CDRL across domains with distinct state and action spaces. The existing approaches can divided into the following main categories:

- (i) Manually designed latent mapping: In (Ammar & Taylor, 2012) and (Ammar et al., 2012), the trajectories are mapped manually and by sparse coding from the source domain and the target domain to a common latent space, respectively. The distance between latent states can then be calculated to find the correspondence of the states from the different domains. In Gupta et al. (2017), the correspondence of the states is found by dynamic time warping and the mapping function which can map the states from two domains to the latent space is found by the correspondence.
- (ii) Learned inter-domain mapping: In the literature (Taylor et al., 2008; Zhang et al., 2021; You et al., 2022; Gui et al., 2023; Zhu et al., 2024), the inter-domain mapping is mainly learned by enforcing dynamics alignment (or termed dynamics cycle consistency in (Zhang et al., 2021)), i.e., aligning the one-step transitions of the two domains. Additional properties have also been incorporated as auxiliary loss functions in learning the inter-domain mapping in the prior works, including domain cycle consistency (Zhang et al., 2021; You et al., 2022), effect cycle consistency (Zhu et al., 2024), maximizing mutual information between states and embeddings (You et al., 2022), and alignment of target-domain rewards with the embeddings (You et al., 2022). Moreover, as the state and action spaces are typically bounded sets and these methods directly map the data samples between the two domains, adversarial learning has been used to restrict the output range of the mapping functions (Zhang et al., 2021; Gui et al., 2023). On the other hand, in (Ammar et al., 2015), the state mapping function is found by Unsupervised Manifold Alignment (Wang & Mahadevan, 2009).

Despite the above progress, the existing approaches all presume that the domains are sufficiently similar and do not have any performance guarantees (and hence can suffer from negative transfer in bad-case scenarios). By contrast, this paper proposes a robust CDRL method that can achieve transfer regardless of source-domain model quality or domain similarity with guarantees.

CDRL across domains with identical state and action spaces. In CDRL, a variety of methods have been proposed for the case where source and target domains share the same state and action spaces but are subject to dynamics mismatch.

- (i) Using the data samples from both source and target domains for policy learning: One popular approach is to use the data from both domains for model updates (Eysenbach et al., 2021; Liu et al., 2022; Xu et al., 2023). For example, for compensating the discrepancy between domains in transition dynamics, (Eysenbach et al., 2021) proposes to modify the reward function, which is learned by an auxiliary domain classifier that distinguishes between the source-domain and target-domain transitions. (Liu et al., 2022) handles the dynamics shift problem in offline RL by augmenting rewards in the source-domain dataset. (Xu et al., 2023) proposes to address dynamics mismatch by a value-guided data filtering scheme, which ensures selective sharing of the source-domain transitions based on the proximity of paired value targets.
- (ii) Explicit domain similarity: (Sreenivasan et al., 2023) proposes to selectively apply direct transfer of the source-domain policy to the target domain based on a learnable similarity metric, which is essentially the TD error of target domain trajectories with source Q function. Moreover, based on the policy invariant explicit shaping (Behboudian et al., 2022), (Sreenivasan et al., 2023) further uses the potential function as a bias term for selecting actions.
- (iii) *Using both Q-functions for the Q-learning updates*: Target Transfer Q-Learning (Wang et al., 2020) calculates the TD error by the source and target domains Q functions in order to select the TD target from the two Q functions.
- (iv) *Domain randomization*: To tackle sim-to-real transfer with dynamics mismatch, domain randomization (Rajeswaran et al., 2016; Peng et al., 2018; Chebotar et al., 2019; Du et al., 2021) and Du et al. (2021) collects data from multiple similar source domains with different

configurations to learn a high-quality policy that can work robustly in a possibly unseen but similar target domain.

D ADDITIONAL EXPERIMENTAL RESULTS

D.1 A TOY EXAMPLES FOR MOTIVATING THE BENEFIT OF CROSS-DOMAIN BELLMAN LOSS

We consider the 3-by-3 grid navigation problem, as shown in Figure 5. In both domains, there are only two actions: 'going top' and 'going right.' The state of the source domain is described in decimal coordinates, while the state of the target domain is described in binary coordinates. The white squares represent obstacles that cannot be traversed. There are three special states: (i) Start state: The episode always begins at this state. (ii) End state: The episode will only end at this state, and the agent will receive an ending reward of +1. (iii) Treasure state: When the agent first navigates to this state, it will receive +0.5 rewards. In other states or at other times

(0, 2)	(1, 2)	(2, 2)		(0, 0, 1, 1)	(0, 1, 1, 1)	(1, 1, 1, 1)
(0, 1)	(1, 1)	(2, 1)		(0, 0, 0, 1)	(0, 1, 0, 1)	(1, 1, 0, 1)
(0, 0)		(2, 0)		(0, 0, 0, 0)		(1, 1, 0, 0)
			'			

(a) Source Domain

(b) Target Domain

Figure 5: Source and target domains of the grid navigation example.

navigating the treasure state, the agent will not receive any reward. In the source domain, the start state, end state, and treasure state are set to (0,0), (0,2), and (2,2), respectively. In the target domain, the start state, end state, and treasure state are set to (0,0,0,0), (0,0,1,1), and (1,1,1,1), respectively. We assume that the source Q-function $Q_{\rm src}$ is optimal in the source domain and the environment discount factor γ is set to 0.99. It is easy to verify that the optimal trajectory of the source domain is $(0,0) \to (0,1) \to (0,2) \to (1,2) \to (2,2)$ and the optimal trajectory of the target domain is $(0,0,0,0) \to (0,0,0,1) \to (0,0,1,1) \to (0,1,1,1) \to (1,1,1,1)$. Consider two trajectories in the source domain: Traj-A, which is the optimal trajectory, and Traj-B, defined as $(0,0) \to (0,1) \to (1,1) \to (1,2) \to (2,2)$. When we map the optimal trajectory of the target domain to Traj-A and the optimal trajectory of the target domain to Traj-B, both mappings result in 0 cycle consistency loss. This suggests that the cycle consistency cannot determine which mapping is superior. This phenomenon results from the unsupervised nature of dynamics cycle consistency. In contrast, when we mapping the optimal trajectory of the target domain to Traj-A yields a cross-domain Bellman-like loss of 0, while mapping the optimal trajectory of the target domain to Traj-B results in a cross-domain Bellman-like loss of 1. Thus, we can achieve optimal mapping results based on the cross-domain Bellman error, while the cycle consistency loss provides sub-optimal mapping results.

D.2 FINAL REWARDS

In this section, we show the asymptotic performance of all baselines and our algorithm. In the experiments, we train all the target-domain models for 500k steps in MuJoCo and 100k steps in Robosuite. The asymptotic performances of all baselines and our algorithm are shown in the following Table 2.

Table 2: Final rewards of QAvatarand all baselines in the experiments.

Algorithm	HalfCheetah	Ant	Door Opening	Table Wiping	Navigation
QAvatar SAC FT PAR CAT-SAC	11586.0 ± 1224.4 10986.0 ± 1821.8 10756.8 ± 1070.8 8097.4 ± 3962.0 8756.5 ± 1264.3 46.1 + 149.9	2858.8 ± 848.0 1620.0 ± 527.2 1644.3 ± 748.2 737.6 ± 45.3 1628.9 ± 200.6 17.1 ± 27.3	216.6 ± 131.3 94.8 ± 23.9 129.9 ± 34.6 33.7 ± 18.6 63.2 ± 33.3 34.7 ± 8.4	76.6 \pm 13.5 47.6 \pm 11.0 42.1 \pm 15.4 17.9 \pm 11.8 23.7 \pm 10.7 55.5 \pm 29.7	38.5 ± 13.2 19.7 ± 13.6 12.5 ± 9.0 0.0 ± 0.0 2.7 ± 2.4 -0.1 ± 0.2
CMD	-253.1 ± 344.1	777.5 ± 144.1	7.8 ± 6.4	0.8 ± 0.4	-0.0 ± 0.0

D.3 ABLATION STUDY: EXPERIMENT RESULT

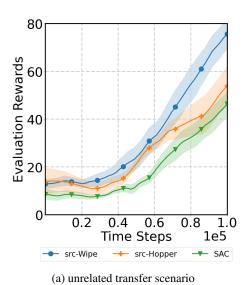
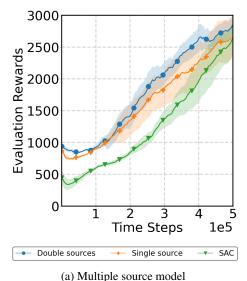


Figure 6: Training curves of unrelated trans-

fer scenario, the source domains are labeled. The target domain is Table-Wiping with robot UR5e.



(a) Multiple source model

Figure 8: Training curves of 2 source domains transfer to target domain.

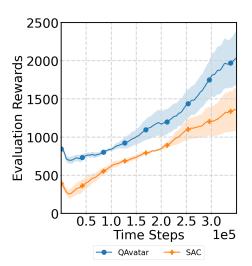
E $\,$ Implementation Details of QAvatar

E.1 PSEUDO CODE OF THE PRACTICAL IMPLEMENTATION OF QAVATAR

In this section, we provide the pseudo code of the practical version of QAvatarin Algorithm 3.

E.2 SOURCE-DOMAIN MODELS AND THEIR PERFORMANCE

For the locomotion tasks including HalfCheetah and Ant, we train each source model for 1M steps. The average performance of the 5 source-domain models (under 5 distinct random seeds) in HalfCheetah and Ant are 7355 ± 2892 and 3689 ± 1013 , respectively. For the Robosuite tasks



(a) non-stationary environment

Figure 7: Training curves in non-stationary Ant-v3.

Algorithm 3 Practical Implementation of QAvatar

- 1: **Require:** Source-domain Q-network Q_{src} , update α frequency N_{α} , batch size N.
- 2: Initialize the state mapping function ϕ , the action mapping function ψ , the initial target-domain policy network $\pi^{(1)}$, entropy coefficient β , replay buffer D, and $\alpha = 0$.
- 3: **for** iteration $t = 1, \dots, T$ **do**

- 4: Interact with the environment and store the transition (s_t, a_t, r_t, s_{t+1}) in the replay buffer D.
- 5: Sample two sets of N transitions, denoted as B_{SAC} and B_{Map} , from the replay buffer D.
- 6: Update the target-domain $\{Q_{tar,1}, Q_{tar,2}\}$ by SAC's critic loss:

$$Q_{\text{tar},j}^{(t)} = \arg\min_{Q_{\text{tar}}} \hat{\mathbb{E}}_{(s,a,r,s') \in B_{\text{SAC}}} \Big[\big| r + \gamma \mathbb{E}_{a' \sim \pi^{(t)}(\cdot|s')} \big[Q_{\text{tar}}(s',a') - \beta \log(\pi(a'|s')) \big] - Q_{\text{tar}}(s,a) \big|^2 \Big].$$
(72)

- 7: Update the state mapping function ϕ and action mapping function ψ by minimizing
- 8: the following loss:

$$\phi^{(t)}, \psi^{(t)} = \arg\min_{\phi, \psi} \hat{\mathbb{E}}_{(s, a, r, s') \in B_{\text{Map}}} \Big[\Big| r + \gamma \mathbb{E}_{a' \sim \pi^{(t)}(\cdot | s')} \big[Q_{\text{src}}(\phi(s'), \psi(a')) \big] - Q_{\text{src}}(\phi(s), \psi(a)) \Big|^2 \Big].$$
(73)

- 9: **if** $t \mod N_{\alpha} = 0$ **then**10: Define $\|\epsilon_{\rm td}^{(t)}\|_D = \hat{\mathbb{E}}_{(s,a,r,s')\in D} \Big[\big| r + \gamma \mathbb{E}_{a'\sim\pi^{(t)}(\cdot|s')} \big[\min_{j=1,2} Q_{{\rm tar},j}^{(t)}(s',a') \big] \min_{j=1,2} Q_{{\rm tar},j}^{(t)}(s,a) \big| \Big],$ 11: $\|\epsilon_{\rm cd}(Q_{\rm src},\phi^{(t)},\psi^{(t)})\|_D = \hat{\mathbb{E}}_{(s,a,r,s')\in D} \Big[\big| r + \gamma \mathbb{E}_{a'\sim\pi^{(t)}(\cdot|s')} \big[Q_{\rm src}(\phi^{(t)}(s'),\psi^{(t)}(a')) \big] Q_{\rm src}(\phi^{(t)}(s),\psi^{(t)}(a)) \big| \Big].$ 12: Update the weight $\alpha = \|\epsilon_{\rm td}^{(t)}\|_D / (\|\epsilon_{\rm cd}(Q_{\rm src},\phi^{(t)},\psi^{(t)})\|_D + \|\epsilon_{\rm td}^{(t)}\|_D).$ 13: **end if**
- 14: Update the target-domain policy π :

$$\pi^{(t+1)} = \arg\min_{\pi} \hat{\mathbb{E}}_{\substack{(s,a,r,s') \in B_{SAC} \\ a' \sim \pi^{(t)}(\cdot|s)}} \left[\beta \log \pi(a'|s) - f^{(t)}(s,a') \right], \tag{74}$$

$$f^{(t)}(s,a') = (1-\alpha) \min_{j=1,2} Q_{\text{tar},j}^{(t)}(s,a') + \alpha Q_{\text{src}}(\phi^{(t)}(s), \psi^{(t)}(a')). \tag{75}$$

15: **end for**

including Door Opening and Table Wiping, we train each source-domain model for 500K steps. The average performance of 5 random seed is 383 ± 139 and 94 ± 16 , respectively. For the navigation environment, we train the model for 500K steps, and the average performance is 39.85.

E.3 INTER-DOMAIN MAPPING NETWORK AUGMENTED WITH A NORMALIZING FLOW MODEL

As discussed in Section 4, a flow-based generative model is employed to transform the outputs of the mapping functions into their corresponding feasible regions. Therefore, there are two architectural paradigms of the flow model can be considered. In the first paradigm, the state and action are concatenated and jointly treated as the codomain of the flow model. This joint formulation is adopted in Cheetah, Ant environment. In the second paradigm, the state and action are modeled separately, with two independent flow models trained respectively for the state and the action. This decoupled formulation is applied in Hopper-v3, Table Wiping, and Door Opening tasks.

F CONFIGURATION DETAILS OF THE EXPERIMENTS

F.1 STATE AND ACTION DIMENSIONS OF BENCHMARK ENVIRONMENTS

We summarize the state and action dimensions of each pair of source-domain and target-domain benchmark tasks in the following Table 3.

Table 3: Dimensions of the source and target domains ("Src" and "Tar" represent the source domain and the target domain.)

Environment	State		Action	
Ziiviioiiiieii	Src	Tar	Src	Tar
HalfCheetah Ant	17 111	23 133	6 8	9 10
Door Opening Table Wiping	46 37	51 34	8 7	7 6
Goal Navigation	40	72	2	12

F.2 MuJoCo and Robosuite Environments

As mentioned in Section 5, We evaluate QAvatarin both MuJoCo and Robosuite environments. In the MuJoCo environments, the source domains of our experiments are the original MuJoCo environments such as HalfCheetah-v3 and Ant-v3. The target domains are the modified MuJoCo environments such as HalfCheetah with three legs and Ant with five legs. In Robosuite environments, We evaluate QAvataron two tasks, including door opening and table wiping. For each task, we consider crossdomain transfer from controlling a Panda robot arm to controlling a UR5e robot arm. These four tasks are illustrated in Figure 9 and 10.

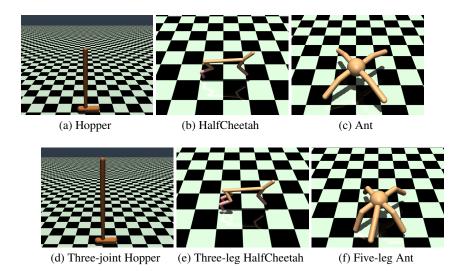


Figure 9: The environments of the source domains and the target domains. (a)-(c): Source domains – Original MuJoCo environments. (d)-(f): Target domains – Modified MuJoCo environments.

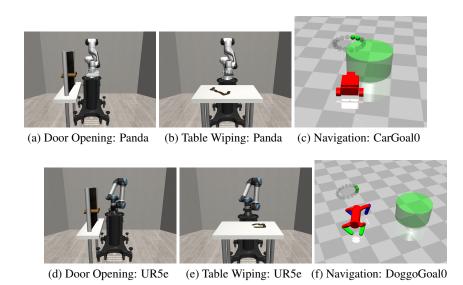


Figure 10: The environments of the source domains and the target domains. (a)-(c): Source domains – Control Panda to solve the tasks in robosuite and Safety-gym CarGoal0. (e)-(h): Target domains – Control UR5e to solve the tasks in robosuite and Safety-gym DoggoGoal0.

F.3 THE IMPLEMENTATION DETAILS OF BASELINES

SAC. The implementation of SAC used in our experiments is released by stable-baselines Raffin et al. (2021). The settings of all hyperparameters except for the discouted factor γ follows the default settings of SAC in the documentation of stable-baselines 3. The discouted factor is set 0.99 in all other MuJoCo environments, which follows the setting shown in Hugging Face. As for in the Robosuite environments, we set the discouted factor to 0.9.

CMD. Since there is no publicly available implementation of CMD, we leverage and adapt the codebase of DCC (Zhang et al., 2021) (https://github.com/sjtuzq/Cycle_Dynamics) and reproduce CMD by following the pseudo code of CMD in its original paper Gui et al. (2023). We follow the setting of the hyperparameters which is revealed in its original paper. Additionally, we change CMD from collecting the fixed amount of data to collecting data continuously for a fair comparison. As for the source model, we use the same model used in our algorithm. Moreover, we observe that the original setting could suffer because the collected trajectories mostly have low returns due to a random behavior policy. Therefore, we consider a stronger version of CMD with target-domain data collected under the target-domain policy, which is induced by the source-domain pre-trained policy and the current inter-domain mappings.

FT. FT can be seen as a standard SAC algorithm with source feature initialization. Specifically, we modify the input and output layers of the source policy to match the target domain's state and action dimensions, using random initialization, while keeping the middle layers with the same weights as the source model. Similarly, for the source Q function, we adjust the input layer to fit the target domain's state and action dimensions with random initialization, while the remaining layers retain the source model's weights. After initialization, we can use SAC algorithm to implement FT.

CAT. We use the authors' implementation (https://github.com/TJU-DRL-LAB/transfer-and-multi-task-reinforcement-learning/tree/main/Single-agent%20Transfer%20RL/Cross-domain%20Transfer/CAT) and use PPO as the target-domain base algorithm following the original paper. For a fair comparison, we use the same source model used in QAvatar. The hyperparameters are shown in the following table and "n epochs" means the number of epochs when optimizing the surrogate loss.

CAT-SAC. As CAT can be integrated with any off-the-shelf RL method, we adapt the original PPO-based CAT to CAT-SAC by using the SAC implementation in Spinning Up Achiam (2018) as the backbone of CAT-SAC. All the SAC-related hyperparameters are the same as those used by SAC and the CAT-related parameters are configured as in the original implementation. For a fair comparison, we use the same source model used by *Q*Avatar.

PAR. We use the authors' implementation (https://github.com/dmksjfl/PAR.git) and consider the offline to online version of PAR, which is more compatible with the CDRL setting in our paper. For the source-domain data required by PAR, we use the samples in the buffer collected during the training of the source-domain policies (shared by QAvatarand other baselines). As a result, to adapt PAR to the more general CDRL setting in our paper, similar to the data pre-processing methods used in handling sequences (Zahavy et al., 2018; Dwarampudi & Reddy, 2019; Morad et al., 2024; ?), we use padding and truncation to handle the differences in state and action dimensions. More specifically,

- **Padding**: If the target domain has n more dimensions than the source, we append n zeros to the end of each source sample.
- **Truncation**: If the target domain has n fewer dimensions than the source, we discard the last n from each source sample.

Note that this design is reasonable, as neither the baselines nor QAvatarhave any knowledge about the physical meaning of each entry in the state or action representations. For the hyperparameters, to ensure a fair comparison with QAvataras well as the baselines CAT-SAC and SAC, we set the ratio between environment interaction and agent training to 1 (i.e., config['tar_env_interact_freq'] in their original code). Other parameters (e.g., beta, weight, etc.) and network architecture follow the recommendations provided in the original PAR paper. In addition, we observe that in some environments, temperature tuning can improve performance. Therefore, we apply temperature tuning during the training process (as adopted by PAR's original code), and select the better one between using and not using temperature tuning as the final result.

Table 4: A list of candidate hyperparameters for Robosuite and MuJoCo.

Parameter	MuJoCo	Robosuite
learning rate	0.0001, 0.0003, 0.0004, 0.0008	0.0001, 0.0003
length of rollouts batch size	500, 2000 50, 100	2000 50, 100, 200
entropy coefficient (ent. coef.)	0.01, 0.002	0.01, 0.002
n epochs num. of hidden layer of encoder/decoder	10, 20 1	5, 10 1
num. of hidden layer of actor/critic	2	2
hidden layer size	256	256

Table 5: Final hyperparameters chosen for each environment.

	learning rate	len. of rollouts	batch size	ent. coef.	n epochs
HalfCheetah	0.0001	500	50	0.002	10
Ant	0.0004	500	50	0.002	10
Robosuite	0.0003	2000	100	0.01	10

F.4 DETAILED CONFIGURATION OF QAVATAR

The base algorithm, SAC, is implemented by stable-baselines 3 Raffin et al. (2021). As for the compute resource, we use NVIDIA GeForce RTX 3090 to do the experiments. The Hyperparameters of QAvatarare shown in the following table. The settings of hyperparameters such as critic/actor learning rate, batch size, buffer size and discounted factor are same as SAC.

Table 6: A list of hyperparameters of QAvatar.

Parameter	Value
critic/actor learning rate	0.0003
state mapping function learning rate	0.01
action mapping function learning rate	0.01
batch size	256
replay buffer size	10^{6}
optimizer	Adam
number of hidden layer of mapping functions	1
hidden layer size	256
update α frequency N_{α}	1000