CROSS-DOMAIN REINFORCEMENT LEARNING UNDER DISTINCT STATE-ACTION SPACES VIA HYBRID Q FUNCTIONS

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

025

026

027

028 029 030

031

Paper under double-blind review

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 domain similarity in RL is not easily identifiable a priori, and hence CDRL can be prone to negative transfer. In this paper, we propose to jointly tackle these two challenges through the lens of hybrid Q functions. Specifically, we propose QAvatar, which combines the Q functions from both the source and target domains with a proper weight decay function. Through this design, we characterize the convergence behavior of QAvatar and thereby 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 extensive experiments, we demonstrate that QAvatar achieves superior transferability across domains on a variety of RL benchmark tasks, such as locomotion and robot arm manipulation, even in the scenarios of potential negative transfer.

1 INTRODUCTION

032 Reinforcement learning (RL) has witnessed significant progress in various challenging domains, such 033 as game playing (Mnih et al., 2015; Silver et al., 2016), robot control (Gu et al., 2017; Kalashnikov 034 et al., 2018), and language models (Ouyang et al., 2022), mainly due to the integration of general RL techniques with advancements in data collection and computation for large-scale training. However, data inefficiency of RL remains one significant obstacle to its deployment in many real-world applications, where online data collection is either costly (e.g., robotics and autonomous driving) or 037 even hazardous (e.g., medical treatments). As one promising solution, cross-domain RL (CDRL) serves as a practical framework to improve the sample efficiency of RL from the perspective of transfer learning, which leverages the data or the pre-trained models from a source domain to enable 040 knowledge transfer to the target domain, under the presumption that the data collection and model 041 training are much less costly in the source domain (e.g., simulators). 042

A plethora of the existing CDRL methods focuses on knowledge transfer across environments that 043 share the same state-action spaces but with different transition dynamics. This setting has been 044 extensively studied from a variety of perspectives, such as domain randomization (Peng et al., 2018), learning similarity metrics (Sreenivasan et al., 2023), reward augmentation (Eysenbach et al., 2021; 046 Liu et al., 2022), and data filtering (Xu et al., 2023). Despite the above progress, to fully realize the 047 promise of CDRL, there are two further fundamental challenges to tackle: (i) Distinct state and/or 048 action spaces between domains: To support flexible transfer across a wide variety of domains, the generic CDRL algorithms are 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 051 to apply direct policy transfer across robot agents of different morphologies (Zhang et al., 2021), which naturally leads to discrepancy in representations. This discrepancy significantly complicates 052 the transfer of either data samples or learned source-domain models. (ii) Unknown domain similarity and negative transfer: Typical CDRL presumes that the source and target domains are sufficiently

054 similar such that effective transfer is achievable. However, in practice, given that the data budget 055 of the target domain is limited, it is rather difficult to determine a priori the similarity of a pair of 056 domains, and this becomes even more challenging when the state-action spaces of the two domains 057 are distinct. Moreover, this issue can also be highlighted by the phenomenon of negative transfer 058 (Weiss et al., 2016; Pan & Yang, 2009), where transfer learning from the source domain can have a negative impact on the target domain. As a consequence, despite that CDRL has been shown to succeed in various scenarios, without a proper design, the performance of CDRL could actually 060 be much worse than the vanilla target-domain model learned without using any source knowledge 061 beyond these good-case scenarios. Notably, to tackle (i), several approaches have been proposed to 062 address such representation discrepancy by learning state-action correspondence, either in the typical 063 RL (You et al., 2022) or unsupervised settings (Zhang et al., 2021; Gui et al., 2023). However, these 064 existing solutions are all oblivious to the issues of domain dissimilarity and negative transfer and 065 therefore do not provide any performance guarantees. As a result, one fundamental research question 066 about CDRL remains largely open: How to achieve efficient and reliable cross-domain transfer in RL 067 across domains of distinct state-action spaces without the knowledge about domain similarity? 068

In this paper, we answer the above question in the affirmative. Specifically, we revisit the cross-069 domain transfer problem in RL from the perspective of mixing the source-domain and target-domain 070 Q functions and propose a new CDRL framework termed QAvatar, where an "avatar", as described 071 in the movie Avatar, refers to a genetically engineered body that is created by combining human 072 DNA with the DNA of the native inhabitants of the alien moon. These avatars allow humans on Earth 073 to remotely control these bodies and quickly adapt to the toxic environment of another planet. By 074 drawing an analogy between the cross-planet transfer of humans and the cross-domain transfer of 075 models in RL, we propose to construct a QAvatar, which updates the target-domain policy based on the weighted combination of the learned target-domain Q function and the given source-domain Q 076 function and learn the state-action correspondence by minimizing a cross-domain Bellman loss. 077

078 To substantiate this idea, we first present a prototypical algorithm of QAvatar in the tabular setting 079 and establish that QAvatar enjoys a nice upper bound on the sub-optimality under a properly designed weight decay function, regardless of the similarity between the source and target domains. This result 081 also suggests that QAvatar can achieve improved sample efficiency of CDRL while preventing the potential negative transfer. Based on these findings, we further propose a practical implementation by integrating the QAvatar algorithm with a neural mapping function based on a normalizing flow 083 model in learning the state-action correspondence. 084

085 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 087 for improving sample efficiency. We then present a prototypical QAvatar algorithm and establish its 088 convergence property, showing that QAvatar can improve sample efficiency while avoiding negative transfer. 2) We further substantiate the QAvatar framework by proposing a practical implementation 089 with a normalizing-flow-based state-action mapping. This further demonstrates the compatibility of 090 QAvatar with off-the-shelf methods for learning state-action correspondence. 3) Through extensive 091 experiments and an ablation study, we show that QAvatar significantly outperforms the benchmark 092 CDRL algorithms in various popular RL benchmark tasks, regardless of the quality of source-domain 093 models and domain similarity. 094

095

RELATED WORK 2

096 097

098 **CDRL** across domains with distinct state and action spaces. The existing approaches can divided 099 into three main categories: (i) Manually designed latent mapping: In (Ammar & Taylor, 2012) and 100 (Ammar et al., 2012), the trajectories are mapped manually and by sparse coding from the source 101 domain and the target domain to a common latent space, respectively. The distance between latent 102 states can then be calculated to find the correspondence of the states from the different domains. 103 In Gupta et al. (2017), the correspondence of the states is found by dynamic time warping and the 104 mapping function which can map the states from two domains to the latent space is found by the 105 correspondence. (ii) Learned inter-domain mapping: In the literature (Taylor et al., 2008; Zhang et al., 2021; Heng et al., 2022; Gui et al., 2023; Zhu et al., 2024), the inter-domain mapping is mainly 106 learned by enforcing dynamics alignment (or termed dynamics cycle consistency in (Zhang et al., 107 2021)), i.e., aligning the one-step transitions of the two domains. Additional properties have also

108 been incorporated as auxiliary loss functions in learning the inter-domain mapping in the prior works, 109 including domain cycle consistency (Zhang et al., 2021; Heng et al., 2022), effect cycle consistency 110 (Zhu et al., 2024), maximizing mutual information between states and embeddings (Heng et al., 111 2022), and alignment of target-domain rewards with the embeddings (Heng et al., 2022). Moreover, 112 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 113 of the mapping functions (Zhang et al., 2021; Gui et al., 2023). On the other hand, in (Ammar 114 et al., 2015), the state mapping function is found by Unsupervised Manifold Alignment (Wang & 115 Mahadevan, 2009). Despite the above progress, the existing approaches all presume that the domains 116 are sufficiently similar and do not have any performance guarantees (and hence can suffer from 117 negative transfer in bad-case scenarios). By contrast, this paper proposes a robust CDRL method that 118 can achieve transfer regardless of source-domain model quality or domain similarity with guarantees. 119

CDRL across domains with identical state and action spaces. In CDRL, a variety of methods 120 have been proposed for the case where source and target domains share the same state and action 121 spaces but are subject to dynamics mismatch. (i) Using the data samples from both source and target 122 domains for policy learning: One popular approach is to use the data from both domains for model 123 updates (Eysenbach et al., 2021; Liu et al., 2022; Xu et al., 2023). For example, for compensating the 124 discrepancy between domains in transition dynamics, (Eysenbach et al., 2021) proposes to modify 125 the reward function, which is learned by an auxiliary domain classifier that distinguishes between the 126 source-domain and target-domain transitions. (Liu et al., 2022) handles the dynamics shift problem 127 in offline RL by augmenting rewards in the source-domain dataset. (Xu et al., 2023) proposes to 128 address dynamics mismatch by a value-guided data filtering scheme, which ensures selective sharing 129 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 130 policy to the target domain based on a learnable similarity metric, which is essentially the TD error 131 of target domain trajectories with source Q function. Moreover, based on the policy invariant explicit 132 shaping (Behboudian et al., 2022), (Sreenivasan et al., 2023) further uses the potential function as 133 a bias term for selecting actions. (iii) Using both Q-functions for the Q-learning updates: Target 134 Transfer Q-Learning (Wang et al., 2020) calculates the TD error by the source and target domains 135 Q functions in order to select the TD target from the two Q functions. (iv) Domain randomization: 136 To tackle sim-to-real transfer with dynamics mismatch, domain randomization (Rajeswaran et al., 137 2016; Peng et al., 2018; Chebotar et al., 2019; Du et al., 2021) and Du et al. (2021) collects data from 138 multiple similar source domains with different configurations to learn a high-quality policy that can 139 work robustly in a possibly unseen but similar target domain.

- 140
- 141 142 143

144 145

3 PRELIMINARIES

146 In this section, we provide the problem formulation and basic building blocks of CDRL as well 147 as the useful notation needed by subsequent sections. For a set \mathcal{X} , we let $\Delta(\mathcal{X})$ denote the set 148 of probability distributions over \mathcal{X} . As in typical RL, we model each environment as an infinite-149 horizon discounted Markov decision process (MDP) denoted by $\mathcal{M} := (\mathcal{S}, \mathcal{A}, P, r, \gamma, \mu)$, where (i) 150 S and A represent the state space and action space, (ii) $P: S \times A \to \Delta(S)$ denotes the transition 151 function, (iii) $r: S \times A \to [-R_{\max}, R_{\max}]$ is the reward function, (iv) $\gamma \in [0, 1)$ is the discounted 152 factor, and (v) $\mu \in \Delta(S \times A)$ denotes the initial state-action distribution. Notably, the use of an 153 initial distribution over states and actions is a standard setting in the literature of natural policy gradient (NPG) (Agarwal et al., 2021; Ding et al., 2020; Yuan et al., 2022; Agarwal et al., 2020; 154 Zhou et al., 2024). Given any policy $\pi : S \to \Delta(A)$, we use $\tau = (s_0, a_0, r_1, \cdots)$ to denote a 155 (random) trajectory generated under π in \mathcal{M} , and the expected total discounted reward under π is 156 defined as $V_{\mathcal{M}}^{\pi}(\mu) := \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | \pi; s_0, a_0 \sim \mu]$. Moreover, as usual, 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 state-action 157 158 visitation distribution (also known as the occupancy measure in the MDP literature) of a policy π as $d^{\pi}(s, a) := (1 - \gamma) (\mu(s, a) + \sum_{t=1}^{\infty} \gamma^t \mathbb{P}(s_t = s, a_t = a; \pi))$, for each (s, a). 159 160

161 **Problem Formulation 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

162 MDP $\mathcal{M}_{tar} := (\mathcal{S}_{tar}, \mathcal{A}_{tar}, P_{tar}, \gamma_{tar}, \gamma, \mu_{tar})^{1}$. Notably, in addition to distinct state and action spaces, 163 the two domains can have different reward functions, transition dynamics, and initial distributions. 164 Here we assume that the two MDPs share the same discounted factor γ , which is rather mild. More-165 over, the trajectories of the two domains are completely unpaired. Let Π_{tar} be the set of all stationary 166 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., $\pi^* := \arg \max_{\pi \in \Pi_{tar}} V^{\pi}_{\mathcal{M}_{tar}}(\mu_{tar})$. To 167 improve sample efficiency via knowledge transfer (compared to learning from scratch), in CDRL, 168 the target-domain agent is granted access to $(\pi_{\rm src}, Q_{\rm src}, V_{\rm src})$, which denotes a policy and the cor-169 responding Q and value functions pre-trained in \mathcal{M}_{src} . Notably, we make no assumption on the 170 quality of $\pi_{\rm src}$ (and hence $\pi_{\rm src}$ may not be optimal to $\mathcal{M}_{\rm src}$), despite that $\pi_{\rm src}$ shall exhibit acceptable 171 performance in practice. 172

In this paper, we focus on designing a reliable CDRL algorithm in the sense 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.

176

177 **Inter-Domain Mapping Functions.** To address the discrepancy in state-action spaces in CDRL, 178 learning an inter-domain mapping function is one common building block of many CDRL algorithms. 179 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 Heng et al. (2022) like cycle-GAN (Zhu et al., 2017), neural networks trained by dynamics alignment of the 181 MDPs (Gui et al., 2023). Moreover, mapping functions have various candidate target spaces, such as a 182 latent space, state or action spaces of the target domain (i.e., from S_{src} , A_{src} to S_{tar} , A_{tar}), and state or 183 action spaces of the source domain (i.e., from S_{tar} , A_{tar} to S_{src} , A_{src}). For example, (Gui et al., 2023) proposed to learn two mapping functions $G_1 : S_{tar} \to S_{src}$ and $G_2 : A_{src} \to A_{tar}$ through dynamics 185 alignment, which infers the unknown mapping between the unpaired trajectories of \mathcal{M}_{src} and \mathcal{M}_{tar} 186 by aligning the one-step state transitions. Specifically, dynamics alignment can be implemented 187 by minimizing the loss function defined as $L(G_1, G_2) = \mathbb{E}_{s_{\text{tar}} \sim \rho, s'_{\text{tar}}, s'_{\text{str}}} [\|s'_{\text{str}} - G_1(s'_{\text{tar}})\|_1]$, where 188 s_{tar} is drawn from some target-domain state distribution ρ and $s'_{\text{tar}} \sim P_{\text{tar}}(\cdot|s_{\text{tar}}, G_2(a_{\text{src}}))$ with 189 $a_{\rm src} \sim \pi_{\rm src}(\cdot | G_1(s_{\rm tar}))$. However, this approach provides no performance guarantee as it can suffer 190 from identification issue due to its unsupervised nature. By contrast, in this work, we propose to learn inter-domain state and action mapping functions in the form of $\phi : S_{tar} \to S_{src}$ and $\psi : A_{tar} \to A_{src}$ 191 by leveraging a cross-domain Bellman-like loss function with guarantees, as described subsequently 192 in Section 4. Moreover, we construct a toy example to show that dynamics cycle consistency could 193 get stuck at a sub-optimal inter-domain mapping while the proposed cross-domain Bellman-like loss 194 can learn a better mapping by considering the target-domain rewards in Appendix C.1. 195

Notation. Throughout this paper, for any real-valued function $h : S \times A \to \mathbb{R}$, for any policy π , 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 any $p \ge 1$, we use $||z||_p$ to denote the ℓ_p -norm of z.

199 200

201 202

203

204

205 206 207

208

4 Methodology

In this section, we first describe the prototypical framework of QAvatar in the tabular setting (i.e., S_{tar} and A_{tar} are finite) and establish convergence guarantees. We then extend this framework to a practical deep RL implementation.

4.1 THE QAVATAR FRAMEWORK

The main idea of QAvatar is to utilize a weighted combination of a learned target-domain Q function and the given source-domain Q function for robust cross-domain knowledge transfer. In this way, QAvatar can enjoy improved sample efficiency in good-case scenarios (e.g., \mathcal{M}_{src} and \mathcal{M}_{tar} are similar) while avoiding potential negative transfer in other scenarios. Specifically, QAvatar consists of the following three major components:

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

Algorithm 1 QAvatar

2: f	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 (2), i.e., $Q_{\text{tar}}^{(t)} \leftarrow \arg \min_{Q_{\text{tar}}} \mathcal{L}_{\text{TD}}(Q_{\text{tar}}; \pi^{(t)}, \mathcal{D})$
5:	Update ϕ and ψ by minimizing (1), i.e., $\phi^{(t)}, \psi^{(t)} \leftarrow \arg \min_{\phi, \psi} \mathcal{L}_{CD}(\phi, \psi; Q_{src}, \pi^{(t)}, \mathcal{D}_{t})$
6:	Update the target-domain policy by adapting NPG to CDRL as in (3).
7: (8:]	and for Return Target-domain policy $\pi_{tar}^{(T)} \sim \text{Uniform}(\{\pi^{(1)}, \cdots, \pi^{(T)}\}).$

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

$$(1)$$

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 (1) looks for a pair of mapping functions ϕ, ψ such that $Q_{\rm src}$ aligns as much with the target-domain transitions as possible. In the special case of $\mathcal{M}_{\rm src} = \mathcal{M}_{\rm tar}$ and ϕ, ψ being identity maps, (1) simply reduces to the standard loss function of temporal difference (TD) learning.

Target-domain Q function: To implement the idea of a hybrid Q function, QAvatar maintains a target-domain Q function Q_{tar}, which is essentially a critic of the current target-domain policy. Specifically, in each iteration t, Q_{tar} is obtained by a policy evaluation step via minimizing the standard TD loss for least-squares policy evaluation (LSPE) (Lagoudakis & Parr, 2001; Yu & Bertsekas, 2009; Lazaric et al., 2012)², i.e.,

$$\mathcal{L}_{\mathrm{TD}}(Q_{\mathrm{tar}};\pi_{\mathrm{tar}},\mathcal{D}_{\mathrm{tar}}) := \hat{\mathbb{E}}_{(s,a,\tau_{\mathrm{tar}},s')\in\mathcal{D}_{\mathrm{tar}}} \Big[\big| r_{\mathrm{tar}} + \gamma \mathbb{E}_{a'\sim\pi_{\mathrm{tar}}} [Q_{\mathrm{tar}}(s',a')] - Q_{\mathrm{tar}}(s,a) \big|^2 \Big], \quad (2)$$

- where $\mathcal{D}_{tar} = \{(s, a, r, s')\}$ denotes target-domain samples.
- NPG-like policy update with a weighted combination of Q functions: The core idea of QAvatar is to leverage both Q_{src} and Q_{tar} to determine policy updates. In the tabular setting, inspired by (Zhou et al., 2024) in the offline-to-online RL literature, we adapt the classic NPG update (Kakade, 2001), which takes an exponential-weight form on the Q function in the policy space (cf. (Agarwal et al., 2021; Xiao, 2022)), to the CDRL setting. In each iteration t,

$$\pi^{(t+1)}(a|s) \propto \pi^{(t)}(a|s) \exp\left(\eta \cdot \left((1 - \alpha(t))Q_{\text{tar}}^{(t)}(s, a) + \alpha(t)Q_{\text{src}}(\phi^{(t)}(s), \psi^{(t)}(a))\right)\right), \quad (3)$$

where $\alpha : \mathbb{N} \to [0, 1]$ is the weight decay function to be configured. Intuitively, $\alpha(t)$ shall be close to one for small t to achieve knowledge transfer from Q_{src} and gradually diminish to zero to escape from potential negative transfer.

The pseudo code of QAvatar is provided in Algorithm 1.

Remark 1. In Line 8 of Algorithm 1, QAvatar outputs the final policy by choosing uniformly at random from the set of all intermediate polices. This is a standard procedure in the optimization literature to connect the average sub-optimality with the performance of output policy. In the experiments, we show that using the last-iterate policy is sufficient and performs well.

4.2 PERFORMANCE GUARANTEES OF QAVATAR

In this section, we formally present the theoretical guarantee of QAvatar and thereby describe how to choose the proper decay parameter $\alpha(\cdot)$. Before stating the theorem, we first describe a useful definition on the coverage in terms of state-action distribution (Zhou et al., 2024).

²These works on LSPE are shown under linear function approximation, which includes the tabular setting as a special case by using one-hot feature vectors.

270 **Definition 1** (Coverage). Given a comparator policy π^{\dagger} in \mathcal{M}_{tar} , we say that π^{\dagger} has coverage $C_{\pi^{\dagger}}$ if 271 for any policy $\pi \in \Pi_{tar}$, we have $\|d^{\pi^{\dagger}}/d^{\pi}\|_{\infty} \leq C_{\pi^{\dagger}}$.

Notably, one can verify that $C_{\pi^{\dagger}}$ is finite if $||d^{\pi^{\dagger}}/\mu_{tar}||_{\infty}$ is finite (given that $||\mu_{tar}/d^{\pi}||_{\infty} \le 1/(1-\gamma)$ for all π , by the definition of d^{π}), and this can be satisfied under an exploratory initial distribution with $\mu_{tar}(s, a) > 0$ for all (s, a), which is one standard assumption in the NPG literature (Agarwal et al., 2021; Ding et al., 2020; Yuan et al., 2022; Agarwal et al., 2020; Zhou et al., 2024). Intuitively, the coverage is needed to enable direct comparison of the Bellman error between policies.

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

279 280
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 defined as $\epsilon_{td}^{(t)}(s, a) := |Q_{tar}^{(t)}(s, a) - r_{tar}(s, a) - \gamma \mathbb{E}_{s' \sim P_{tar}(\cdot | s, a), a' \sim \pi^{(t)}(\cdot | s')} [Q_{tar}^{(t)}(s', a')]|.$

Definition 3 (Cross-Domain Bellman Error). Given a source-domain Q_{src} , for each state-action pair (s, a) and $t \in \mathbb{N}$, the cross-domain Bellman error $\epsilon_{src,be}^{(t)}(s, a; Q_{src})$ is defined as $\epsilon_{src,be}^{(t)}(s, a; Q_{src}) := [Q_{src}(\phi^{(t)}(s), \psi^{(t)}(a)) - r_{tar}(s, a) - \gamma \mathbb{E}_{s' \sim P_{tar}(\cdot|s, a), a' \sim \pi^{(t)}(\cdot|s')} [Q_{src}(\phi^{(t)}(s'), \psi^{(t)}(a'))]].$

Definition 4 (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 $f^{(t)}(s, a) := (1 - \alpha(t))Q_{tar}^{(t)}(s, a) + \alpha(t)Q_{src}(s, a)$, where $\alpha : \mathbb{N} \to [0, 1]$ is the weight decay function.

Below we use $\|\epsilon_{\text{src, be}}^{(t)}(Q_{\text{src}})\|_{\infty}$ and $\|\epsilon_{\text{td}}^{(t)}\|_{\infty}$ as shorthand for $\|\epsilon_{\text{src, be}}^{(t)}(\cdot, \cdot; Q_{\text{src}})\|_{\infty}$ and $\|\epsilon_{\text{td}}^{(t)}(\cdot, \cdot)\|_{\infty}$, and we use $\mu_{\text{tar,min}}$ as a shorthand for $\min_{s,a} \mu_{\text{tar}}(s, a)$. We are ready to present the main theoretical result, and the detailed proof is provided in Appendix B.

Theorem 1. (Average Sub-Optimality) Under the QAvatar in Algorithm 1 and Assumption 1, given any fixed learning rate $\eta > 0$, the average sub-optimality over T iterations can be upper bounded as

$$\frac{1}{T} \sum_{t=1}^{T} \left(V^{\pi^{*}}(\mu_{tar}) - V^{\pi^{(t)}}(\mu_{tar}) \right) \leq \underbrace{\frac{1}{(1-\gamma)T} \sum_{t=1}^{T} \mathbb{E}_{(s,a)\sim d^{\pi^{*}}} \left[\max_{a'} \bar{f}^{(t)}(s,a') \right] + \frac{\log |\mathcal{A}_{tar}|}{(1-\gamma)T\eta}}_{(a)}}_{(a)} + \underbrace{\frac{C_{0}}{T} \sum_{t=1}^{T} \alpha(t) \|\epsilon_{src, \ be}^{(t)}(Q_{src})\|_{\infty}}_{(b)} + \underbrace{\frac{C_{0}}{T} \sum_{t=1}^{T} (1-\alpha(t)) \|\epsilon_{td}^{(t)}\|_{\infty}}_{(c)}}_{(c)}, \quad (4)$$

where $C_0 := 2\sqrt{C_{\pi^*}}/((1-\gamma)^2 \mu_{tar, min})$ and $\bar{f}^{(t)}(s, a) := f^{(t)}(s, a) - f^{(t)}(s, \pi^{(t)}(s))$.

Notably, in (4), the term (a) reflects the learning progress of NPG, the (b) reflects the effect of cross-domain transfer, and (c) indicates the error of policy evaluation for the target-domain policy. The term (c) reflects the sample complexity of the standard least-squares TD-based policy evaluation (Lagoudakis & Parr, 2001; Lagoudakis et al., 2002; Yu & Bertsekas, 2009; Lazaric et al., 2012) and can be made small with sufficient samples (i.e., sufficiently large $N_{tar}^{(t)}$).

Key Implications of Theorem 1.

286

287

288 289

290 291

292

293

305 306

307

308

309

310

311

312

313 • Positive transfer indeed reduces the upper bound of average sub-optimality: For didactic 314 purposes, consider an ideal case in the sense that $Q_{\rm src}$ is optimal in the source domain and there always exists a perfect inter-domain mapping ϕ^* and ψ^* such that $L_{CD}(\phi^*, \psi^*; Q_{src}, \pi_{tar}, \mathcal{D}_{tar}) = 0$ 315 under any policy π_{tar} . In this case, the positive transfer perfectly happens. Let $\alpha(t)$ be close to 316 one initially in the first T iterations and let η be sufficiently large. We can observe that term (b) 317 in (4) is always zero, since $\epsilon_{\rm src, be}(Q_{\rm src})$ is always zero. Regarding the term (a) in (4), since $\alpha(t)$ 318 is initially close to one, we have $\overline{f^{(t)}(s,a)} \approx Q_{\text{src}}(\phi^{(t)}(s),\phi^{(t)}(a)) - Q_{\text{src}}(\phi^{(t)}(s),\phi^{(t)}(\pi^{(t)}(s)))$. By the policy update rule in (3), optimality of source Q function Q_{src} , and the fact that perfect 319 320 inter-domain mapping exists, we know that $\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{(s,a) \sim d^{\pi^*}} \left[\max_{a'} \bar{f}^{(t)}(s,a') \right]$ must be small 321 since $\pi(s)$ would be quickly updated to close to $\pi^*(s)$ for any $s \in S_{tar}$. Moreover, the term (c) in 322 (4) shall also be small initially since $\alpha(t)$ is close 1. By combining all the above, we can conclude 323 that the bound on sub-optimality gap is small in the positive transfer regime.

• QAvatar can avoid getting stuck in the negative transfer regime: Consider a negative transfer case, where $\mathcal{L}_{CD}(\phi, \psi; Q_{src}, \pi_{tar}, \mathcal{D}_{tar})$ is always large under any policy π_{tar} and inter-domain mappings ϕ and ψ . As a result, $\|\epsilon_{\text{src, be}}^{(t)}(Q_{\text{src}})\|_{\infty}$ is large. In this case, given that $\alpha(t)$ is a decreasing function, $\alpha(t)$ shall be close to 0 under large t. We can observe that the term (b) in (4) is small (even if $\|\epsilon_{\text{src, be}}^{(t)}(Q_{\text{src}})\|_{\infty}$ remains large). Note that the policy update rule in (3) reduces to the original NPG based on the target-domain critic since $\alpha(t)$ is close to 0. For the term (a) in (4), we have $\bar{f}^{(t)}(s,a) \approx Q^{(t)}(s,a) - Q^{(t)}(s,\pi^{(t)}(s))$ and the term (a) in (4) reduces to the standard bound for NPG. Similarly, the term (c) in (4) reduces to the standard TD error since $\alpha(t)$ is close to 0. By combining all the above, we conclude that under QAvatar , the bound on sub-optimality gap would not be continuously dominated by the term (b) in (4) in the negative transfer regime.

Remark 2. Note that the proof of Theorem 1 bears some high-level resemblance with (Zhou et al., 2024) as they also use NPG in their hybrid actor-critic (HAC) algorithm. That said, QAvatar is fundamentally different from HAC in two aspects: (i) QAvatar addresses cross-domain transfer while HAC focuses on using offline and online data from the same domain. (ii) QAvatar utilizes the hybrid Q function while HAC applies a hybrid squared error regression loss (i.e., the sum of TD errors calculated from both offline and online data).

341
3424.3PRACTICAL IMPLEMENTATION OF QAVATAR

We extend the *Q*Avatar framework in Algorithm 1 to a practical deep RL implementation for continuous state and action spaces by applying the following design choices. The pseudo code is provided in Algorithm 2 in Appendix.

• Learning the target-domain policy and the Q function. To go beyond the tabular setting and handle continuous state and action spaces, we extend QAvatar by first connecting NPG with soft policy iteration (SPI) (Haarnoja et al., 2018). In the entropy-regularized RL setting, SPI has been shown 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 function with the weighted combination of Q_{tar} and Q_{src} of QAvatar . Regarding the weight decay function $\alpha(t)$, based on the theoretical result, we set $\alpha(t) = t^{-\beta}$ with $\beta > 0$ in the experiments.

• Learning the inter-domain mapping functions with an augmented flow model. Similar to the tabular setting, we learn the inter-domain mappings by minimizing the cross-domain Bellman loss. Notably, in practical RL problems, the state and action spaces are mostly bounded sets. As a result, we need to ensure that the outputs of the inter-domain mappings $\phi : S_{tar} \rightarrow S_{src}$ and $\psi : A_{tar} \rightarrow A_{src}$ fall within the feasible regions. As mentioned in Section 2, adversarial learning is widely adopted to solve this practical problem in the existing literature (Taylor et al., 2008; Zhang et al., 2021; Gui et al., 2023; Zhu et al., 2024). However, we observe that adversarial learning could suffer from unstable training process in practice. Therefore, we use the method proposed by (Brahmanage et al., 2024) and train a normalizing flow model that can map the outputs of the mapping functions to the feasible regions.

5 EXPERIMENTS

In this section, we show that QAvatar achieves effective cross-domain transfer and improves the sample efficiency on various RL benchmark tasks. Moreover, we demonstrate that QAvatar can still perform well even with the existence of negative transfer between the source and target domains.
 Unless stated otherwise, all the results reported in this section are averaged over 5 random seeds.

5.1 EXPERIMENTAL SETTINGS

Benchmark CDRL Methods. We compare the performance of *Q*Avatar with various recent CDRL
 benchmark algorithms under distinct state-action spaces, including Dynamics Cycle-Consistency
 (DCC) (Zhang et al., 2021), Cross-Morphology-Domain Policy Adaptation (CMD) (Gui et al., 2023),
 and Cross-domain Adaptive Transfer (CAT) (Heng et al., 2022). For a fair comparison, all these

378 benchmark methods and QAvatar use the same set of source-domain models (i.e., the policy and the 379 corresponding Q-networks), which are pre-trained by using SAC in the source domain. However, 380 the original DCC is implemented in a batch setting, i.e., a fixed number of trajectories are collected 381 for learning the inter-domain mappings. For a fair comparison, we adapt the DCC in an online 382 setting, i.e., learning while iteratively collecting new trajectories. Regarding CMD, we observe that the original setting could suffer because the collected trajectories mostly have low returns due to a 383 random behavior policy. Therefore, we consider a stronger version of CMD with target-domain data 384 collected under the target-domain policy, which is induced by the source-domain pre-trained policy 385 and the current inter-domain mappings. 386

Moreover, to demonstrate the sample efficiency, we also compare *Q*Avatar with the standard SAC (Haarnoja et al., 2018), which learns from scratch in the target domain, as well as with the direct Fine-Tuning (FT) upon the source models (Ha et al., 2024), which can be viewed as using the standard SAC with source feature initialization. Both methods can serve as reasonably competitive baselines. The hyperparameters are provided in Appendix E.

- Evaluation Environments. We evaluate QAvatar in
 three types of RL benchmark environments:
- 394 • Locomotion: We use the standard MuJoCo environ-395 ments, including Hopper-v3, HalfCheetah-v3 and Ant-396 v3, as the source domains and follow the same pro-397 cedure as in (Zhang et al., 2021; Xu et al., 2023) to modify them for the target domains. Moreover, we 399 consider the Centipede environments in CAT (Heng 400 et al., 2022), using CentipedeFour as the source do-401 main and CentipedeSix as the target domain. Addition-402 ally, for evaluation in the scenario of negative transfer, 403 we use Inverted Pendulum-v2 and the two-joint robot 404 arm Reacher-v2 as the source domains and employ Inverted Double Pendulum-v2 and the three-joint robot 405 arm Reacher-v2 as the respective target domains. The 406 details about the morphology is in Appendix E. 407

Table 1: Dimensionalities of the source
and target domains ("Src" and "Tar" rep-
resent the source domain and the target
domain.

Environment	Sta	ate	Action	
	Src	Tar	Src	Tar
Hopper	11	13	3	4
HalfCheetah	17	23	6	9
Ant	111	133	8	10
Centipede	97	139	10	16
IP / Modified IDP	4	11	1	1
Reacher	11	14	2	2
Block Lifting	42	47	8	7
Door Opening	46	51	8	7
Table Wiping	37	34	7	6
Merging / Highway	12	16	2	2

- **Robot arm manipulation**: We use the environments
 provided by Robosuite, a popular package for robot
- learning released by (Zhu et al., 2020). We evaluate our algorithm on three tasks, including block
 lifting, 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.
- Autonomous driving: We use the environments provided by BARK-ML (Bernhard et al., 2020). Notably, these environments are highly unpredictable due to the complex traffic situations and driver behaviors encoded in the BARK-ML behavior model. For cross-domain transfer, we use merging-v0 as the source domain and highway-v0 as the target domain.
 - Table 1 provides the dimensionalities of the state and action spaces in all the tasks.

Reproduction and Sanity Checks for DCC, CMD, and CAT. Regarding CAT and DCC, we directly use the official implementation provided by the original papers. Moreover, as there is no CMD implementation available, we reproduce CMD by referring to the source code of DCC as these two algorithms are similar. As a sanity check, we evaluate DCC and CMD in multiple MuJoCo tasks and confirm that our reproduced scores are indeed close to those reported in the original papers (despite that DCC and CMD do not perform well due to their unsupervised nature). Similarly, a sanity check for CAT on the Centipede task confirms the reproduced score. The details are in Appendix C.2.

426 427 428

418

5.2 EXPERIMENTAL RESULTS

Does QAvatar improve data efficiency? As shown by the training curves in Figure 1, we observe that QAvatar achieves improved data efficiency via cross-domain transfer than SAC throughout the training process in all the MuJoCo, Robosuite, and BARK-ML tasks, despite that these tasks have rather different dimensionalities as shown in Table 1. CAT achieves moderate performance

456

457 458 459



Figure 1: The training curves of *Q*Avatar and the benchmark methods: (a)-(d): Locomotion tasks in MuJoCo; (e)-(g): Robot arm manipulation tasks in Robosuite; (h) Autonomous driving tasks in BARK-ML.

in Table Wiping, Centipede, and Highway but does not learn effectively in the other tasks. These
 appear reasonable as CAT has no performance guarantees and can suffer if the source and target
 are rather dissimilar, despite that CAT applies policy gradient with target-domain rewards to align
 the inter-domain mapping with the target domain. On the other hand, FT typically achieves slight
 improvement in data efficiency than SAC in MuJoCo but slower learning progress in Robosuite. We
 conjecture that this is because distinct robot arms in Robosuite lead to more dissimilar state-action
 representations and hence require more fine-tuning steps.

467 Regarding CMD, it cannot obtain good returns in most of the tasks. Notably, in some environments 468 like Ant, CMD appears very unstable due to its adversarial learning module for restricting the output 469 of their mapping functions. DCC does not perform well in most tasks, including HalfCheetah. This 470 trend is similar to that in the original paper (Zhang et al., 2021). The rewards obtained by DCC in our experiments are slightly lower than those shown in (Zhang et al., 2021) despite that we try our 471 best to reproduce their results. To strengthen our argument, we offer a comparison of the original 472 and reproduced results in Appendix C.2. We conjecture that the undesired performance of CMD and 473 DCC results from that they learn in an unsupervised manner and hence does not take target-domain 474 rewards into account. 475

Additionally, when we consider the time to threshold metric, our algorithm requires 298k fewer
steps to achieve the threshold than SAC does in the best case. When we consider the asymptotic
performance metric, our algorithm can obtain higher final rewards than SAC. The results of these two
metrics are shown in the Appendix C.3 and C.4.

480 Does QAvatar still perform reliably well when negative transfer is likely to happen? We 481 construct negative transfer scenarios by modifying the environment configurations via swapping 482 action encodings, as described below: (i) In the Reacher environment, we use a two-joint robot arm 483 as the source domain and a three-joint arm as the target domain. To match the action dimensions, the 484 middle joint of the three-joint arm is disabled (i.e., its action is always set to 0). We then alter the 485 three-joint arm's configuration by swapping the encoding of "clockwise" and "counterclockwise" 486 actions (termed "Modified Reacher"). (ii) Similarly, we use Inverted Pendulum (IP) as the source domain and Inverted Double Pendulum (IDP) as the target domain. Then, we modify the configuration of IDP by swapping the encodings of the actions "left" and "right" (termed "Modified IDP").

As a result, negative transfer shall easily 489 occur if we deactivate the inter-domain 490 action mapping ψ of QAvatar such that 491 QAvatar cannot learn by simply mapping 492 the "clockwise" in three-joint Reacher to 493 "counterclockwise" in two-jointed Reacher 494 and "left" in IDP to "right" in IP. As shown 495 in Figure 2, we observe that QAvatar out-496 performs all CDRL benchmarks algorithms and exhibits a learning curve similar to 497 SAC and FT, despite the negative transfer 498 scenarios. This confirms that QAvatar can 499 indeed perform reliably due to the use of 500 hybrid Q function. 501



502 Is QAvatar sensitive to the decay func-503 tion? We evaluate QAvatar with $\alpha(t)$ as

511 512 513

514

515

516

517

518

519

521

522

523 524

531 532 Figure 2: The training curves of *Q*Avatar and the benchmark methods in the negative transfer cases.

504 $1/\sqrt{t}$, 1/t, and $1/t^{1.5}$. As shown in Figure 3, QAvatar can learn successfully regardless of the choice 505 of $\alpha(t)$ and appear consistently favorable under all these choices of $\alpha(t)$.

Does QAvatar still perform reliably with a source-domain model of lower quality? We further
 run QAvatar with low-quality source-domain Q networks, which are pre-trained only for five thousand
 steps in both Hopper and Door Opening. As shown by Figure 4, we find that despite QAvatar is
 affected by the low-quality source model initially, it can quickly catch up and achieve total reward
 comparable to SAC. This appears consistent with the theoretical result in Theorem 1.



Figure 3: The training curves of QAvatar under different decay functions α .

Figure 4: The training curves of *Q*Avatar with a high-quality and a low-quality source model.

6 CONCLUDING REMARKS AND LIMITATIONS

In this paper, we present QAvatar, the first CDRL method that can handle distinct state-action representations between domains with performance guarantees. Based on the idea of combining the source-domain and target-domain Q functions, QAvatar achieves robust knowledge transfer and tackles the negative transfer issue. Through extensive experiments, we show that QAvatar indeed serves as a promising and generic solution to cross-domain transfer in RL. One limitation of this work is that we follow the standard CDRL formulation and consider only one source domain and one target domain. Extending the idea of QAvatar to achieve knowledge transfer from multiple source and target domains is a promising future research direction.

540 BIBLIOGRAPHY

548

549

- 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. *Advances in Neural Information Processing Systems*, 33: 13399–13412, 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*, 22(98):1–76, 2021.
- Haitham B Ammar, Karl Tuyls, Matthew E Taylor, Kurt Driessens, and Gerhard Weiss. Reinforcement
 learning transfer via sparse coding. In *Proceedings of the International Conference on Autonomous Agents and Multiagent Systems*, volume 1, pp. 383–390, 2012.
- Haitham Bou Ammar and Matthew E Taylor. Reinforcement learning transfer via common subspaces.
 In Adaptive and Learning Agents: International Workshop, ALA 2011, Held at AAMAS, pp. 21–36.
 Springer, 2012.
- Haitham Bou Ammar, Eric Eaton, Paul Ruvolo, and Matthew Taylor. Unsupervised cross-domain transfer in policy gradient reinforcement learning via manifold alignment. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 29, 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, pp. 1–14, 2022.
- Julian Bernhard, Klemens Esterle, Patrick Hart, and Tobias Kessler. BARK: Open behavior bench marking in multi-agent environments. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 6201–6208, 2020.
- Janaka Brahmanage, Jiajing Ling, and Akshat Kumar. FlowPG: Action-constrained Policy Gradient with Normalizing Flows. *Advances in Neural Information Processing Systems*, 36, 2024.
- Shicong Cen, Chen Cheng, Yuxin Chen, Yuting Wei, and Yuejie Chi. Fast global convergence of natural policy gradient methods with entropy regularization. *Operations Research*, 70(4): 2563–2578, 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 (ICRA)*, pp. 8973–8979, 2019.
- Dongsheng Ding, Kaiqing Zhang, Tamer Basar, and Mihailo Jovanovic. Natural policy gradient
 primal-dual method for constrained markov decision processes. *Advances in Neural Information Processing Systems*, 33:8378–8390, 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 (ICRA)*, pp. 1290–1296, 2021.
- 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.
- Shixiang Gu, Ethan Holly, Timothy Lillicrap, and Sergey Levine. Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. In *IEEE International Conference on Robotics and Automation (ICRA)*, pp. 3389–3396, 2017.
- ⁵⁹³ 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*, 167:104–117, 2023.

594 595 596	Abhishek Gupta, Coline Devin, YuXuan Liu, Pieter Abbeel, and Sergey Levine. Learning invariant feature spaces to transfer skills with reinforcement learning. In <i>International Conference on Learning Representations</i> , 2017.
597 598 599 600	Seokhyeon Ha, Sunbeom Jeong, and Jungwoo Lee. Domain-aware fine-tuning: Enhancing neural network adaptability. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pp. 12261–12269, 2024.
601 602 603	Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In <i>International Conference on Machine Learning</i> , pp. 1861–1870, 2018.
604 605 606 607	You Heng, Tianpei Yang, Yan Zheng, Jianye Hao, and Matthew E. Taylor. Cross-domain adaptive transfer reinforcement learning based on state-action correspondence. In <i>Conference on Uncertainty in Artificial Intelligence</i> , 2022.
608 609	Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In <i>International Conference on Machine Learning</i> , pp. 267–274, 2002.
610 611 612	Sham M Kakade. A natural policy gradient. Advances in Neural Information Processing Systems, 14, 2001.
613 614 615 616	Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang, Deirdre Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, et al. Scalable deep reinforcement learning for vision-based robotic manipulation. In <i>Conference on Robot Learning</i> , pp. 651–673, 2018.
617 618	Michail G Lagoudakis and Ronald Parr. Model-free least-squares policy iteration. Advances in Neural Information Processing Systems, 14, 2001.
619 620 621 622	Michail G Lagoudakis, Ronald Parr, and Michael L Littman. Least-squares methods in reinforcement learning for control. In <i>Methods and Applications of Artificial Intelligence: Second Hellenic Conference on AI</i> , pp. 249–260. Springer, 2002.
623 624	Alessandro Lazaric, Mohammad Ghavamzadeh, and Rémi Munos. Finite-sample analysis of least-squares policy iteration. <i>Journal of Machine Learning Research</i> , 13:3041–3074, 2012.
625 626 627	Jinxin Liu, Zhang Hongyin, and Donglin Wang. DARA: Dynamics-Aware Reward Augmentation in Offline Reinforcement Learning. In <i>International Conference on Learning Representations</i> , 2022.
628 629 630	Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. <i>Nature</i> , 518(7540):529–533, 2015.
631 632 633 634	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. <i>Advances in Neural Information Processing Systems</i> , 35: 27730–27744, 2022.
635 636 637	Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 22(10):1345–1359, 2009.
638 639 640	Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Sim-to-real transfer of robotic control with dynamics randomization. In <i>IEEE International Conference on Robotics and Automation (ICRA)</i> , pp. 3803–3810, 2018.
641 642 643 644 645	Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dormann. Stable-baselines3: Reliable reinforcement learning implementations. <i>Journal of Machine Learning Research</i> , 22(268):1–8, 2021. URL http://jmlr.org/papers/v22/20-1364.html.
646 647	Aravind Rajeswaran, Sarvjeet Ghotra, Balaraman Ravindran, and Sergey Levine. Epopt: Learning robust neural network policies using model ensembles. In <i>International Conference on Learning Representations</i> , 2016.

648 649 650 651	David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of Go with deep neural networks and tree search. <i>Nature</i> , 529(7587):484–489, 2016.
652 653 654 655	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 <i>PRL Workshop Series –Bridging the Gap Between AI Planning and Reinforcement Learning</i> , 2023.
656 657 658 659	Matthew E Taylor, Gregory Kuhlmann, and Peter Stone. Autonomous transfer for reinforcement learning. In <i>Proceedings of International Joint Conference on Autonomous Agents and Multiagent Systems (AAMAS)</i> , pp. 283–290, 2008.
660 661	Chang Wang and Sridhar Mahadevan. Manifold alignment without correspondence. In <i>IJCAI</i> , volume 2, pp. 3, 2009.
662 663 664	Yue Wang, Yuting Liu, Wei Chen, Zhi-Ming Ma, and Tie-Yan Liu. Target transfer Q-learning and its convergence analysis. <i>Neurocomputing</i> , 392:11–22, 2020.
665 666 667	Karl Weiss, Taghi M Khoshgoftaar, and DingDing Wang. A survey of transfer learning. <i>Journal of Big data</i> , 3:1–40, 2016.
668 669 670	Lin Xiao. On the convergence rates of policy gradient methods. <i>Journal of Machine Learning Research</i> , 23(282):1–36, 2022.
671 672 673 674	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. <i>Advances in Neural Information</i> <i>Processing Systems</i> , 36, 2023.
675 676 677 678	Heng You, Tianpei Yang, Yan Zheng, Jianye Hao, and Matthew E Taylor. Cross-domain adaptive transfer reinforcement learning based on state-action correspondence. In <i>Uncertainty in Artificial Intelligence</i> , pp. 2299–2309, 2022.
679 680 681	Huizhen Yu and Dimitri P Bertsekas. Convergence results for some temporal difference methods based on least squares. <i>IEEE Transactions on Automatic Control</i> , 54(7):1515–1531, 2009.
682 683 684	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 <i>International Conference on Learning Representations</i> , 2022.
685 686 687 688	Qiang Zhang, Tete Xiao, Alexei A Efros, Lerrel Pinto, and Xiaolong Wang. Learning cross-domain correspondence for control with dynamics cycle-consistency. In <i>International Conference on Learning Representations</i> , 2021.
689 690 691	Yifei Zhou, Ayush Sekhari, Yuda Song, and Wen Sun. Offline data enhanced on-policy policy gradient with provable guarantees. In <i>International Conference on Learning Representations</i> , 2024.
692 693 694 695	Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In <i>Proceedings of the IEEE International Conference on Computer Vision</i> , pp. 2223–2232, 2017.
696 697 698	Ruiqi Zhu, Tianhong Dai, and Oya Celiktutan. Cross domain policy transfer with effect cycle- consistency. In <i>IEEE International Conference on Robotics and Automation (ICRA)</i> , 2024.
699 700 701	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. <i>arXiv preprint arXiv:2009.12293</i> , 2020.

A SUPPORTING LEMMAS

Lemma 1 (Performance difference lemma). For any two policies π and π' , for any state s, 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).

Lemma 2. Suppose $f^{(t)}$ and $\pi^{(t)}$ denote the cross-domain value functions and the policies at iteration *t*. Then, for any learning rate η and policy π^* , we have

$$\sum_{t=1}^{T} \mathbb{E}_{(s,a)\sim d^{\pi^*}} \left[f^{(t)}(s,a) - f^{(t)}(s,\pi^{(t)}(s)) \right]$$

$$\leq \sum_{t=1}^{T} \mathbb{E}_{(s,a)\sim d^{\pi^*}} \left[\max_{a'} f^{(t)}(s,a') - f^{(t)}(s,\pi^{(t)}(s)) \right] + \frac{\log |\mathcal{A}_{tar}|}{\eta}.$$

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

$$\pi^{(t+1)}(a \mid s) = \frac{\pi^{(t)}(a \mid s) \exp\left(\eta \bar{f}^{(t)}(s, a)\right)}{\sum_{a'} \pi^{(t)}\left(a' \mid s\right) \exp\left(\eta \bar{f}^{(t)}\left(s, a'\right)\right)}.$$
(5)

Let $Z_t = \sum_{a'} \pi^{(t)} (a' \mid s) \exp(\eta \bar{f}^{(t)}(s, a'))$. By multiplying both sides of (5) by Z_t taking the logarithm, and then taking the expectation on both sides w.r.t $(s, a) \sim d^{\pi^*}$, we obtain

$$\mathbb{E}_{(s,a)\sim d^{\pi^*}}\left[\eta\bar{f}^{(t)}(s,a)\right] = \mathbb{E}_{(s,a)\sim d^{\pi^*}}\left[\log Z_t + \log \pi^{(t+1)}(a\mid s) - \log \pi^{(t)}(a\mid s)\right].$$
 (6)

Next, we bound the term $\log Z_t$. Since the $\log(\cdot)$ is an increasing function, we have

$$\log Z_{t} = \log \left(\sum_{a' \in \mathcal{A}} \pi \left(a' \mid s \right) \exp \left(\eta^{(t)} \bar{f}^{(t)} \left(s, a' \right) \right) \right)$$
$$\leq \log \left(\max_{a' \in \mathcal{A}} \exp \left(\eta \bar{f}^{(t)} \left(s, a' \right) \right) \right)$$
$$\leq \max_{a' \in \mathcal{A}} \left(\eta \bar{f}^{(t)} \left(s, a' \right) \right) = \eta \max_{a' \in \mathcal{A}} \bar{f}^{(t)} \left(s, a' \right).$$

Then, we have

$$\mathbb{E}_{(s,a)\sim d^{\pi^*}}\left[\eta \bar{f}^{(t)}(s,a)\right] \leq \mathbb{E}_{(s,a)\sim d^{\pi^*}}\left[\log \pi^{(t+1)}(a\mid s) - \log \pi^{(t)}(a\mid s) + \eta \max_{a'} \bar{f}^{(t)}(s,a')\right].$$
(7)

By taking the summation over iterations on both side of (7),

$$\sum_{t=1}^{T} \mathbb{E}_{(s,a)\sim d^{*}} \left[\eta \bar{f}^{(t)}(s,a) \right]$$

$$\leq \sum_{t=1}^{T} \mathbb{E}_{(s,a)\sim d^{\pi^{*}}} \left[\eta \max_{a'} \bar{f}^{(t)}(s,a') \right] + \mathbb{E}_{(s,a)\sim d^{\pi^{*}}} \left[\log \pi^{(T+1)}(a \mid s) - \log \pi^{(1)}(a \mid s) \right].$$

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

753
754
755
$$\sum_{t=1}^{T} \mathbb{E}_{(s,a)\sim d^{\pi^*}} \left[\bar{f}^{(t)}(s,a) \right] \leq \sum_{t=1}^{T} \mathbb{E}_{(s,a)\sim d^{\pi^*}} \left[\max_{a'} \bar{f}^{(t)}(s,a') \right] + \frac{\log |\mathcal{A}_{tar}|}{\eta}.$$

Lemma 3 ((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 $f : S \times \mathcal{A} \rightarrow \mathbb{R}$, we have

$$\mathbb{E}_{\tau}\left[\sum_{t=0}^{\infty}\gamma^{t}f(s_{t},a_{t})\right] = \frac{1}{1-\gamma}\mathbb{E}_{(s,a)\sim d^{\pi}}\left[f(s,a)\right].$$
(8)

Lemma 4 (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

$$\frac{p_k(s,a)}{d^{\pi}(s,a)} \le \frac{1}{(1-\gamma)\mu(s,a)}.$$
(9)

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

$$d^{\pi}(s,a) := (1-\gamma) \Big(\mu(s,a) + \sum_{t=1}^{\infty} \gamma^t P(s_t = s, a_t = a; \pi) \Big) \equiv \sum_{t=0}^{\infty} \gamma^t P(s_t = s, a_t = a; \pi).$$
(10)

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

$$p_k(s,a) = \sum_{(s_0,a_0)} \mathbb{P}(s_{\text{next},k} = s, a_{\text{next},k} = a | s_0, a_0; \pi) d^{\pi}(s_0, a_0)$$
(11)

$$= \sum_{(s_0, a_0)} \mathbb{P}(s_{\text{next}, k} = s, a_{\text{next}, k} = a | s_0, a_0; \pi) \cdot (1 - \gamma) \cdot \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s_0, a_t = a_0; \pi)$$
(12)

$$= (1 - \gamma) \cdot \sum_{t=0}^{\infty} \gamma^t \sum_{s_0, a_0} \mathbb{P}(s_{\text{next}, k} = s, a_{\text{next}, k} = a | s_0, a_0; \pi) \cdot \mathbb{P}(s_t = s_0, a_t = a_0; \pi)$$
(13)

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

$$(14)$$

Then, we have

$$\frac{p_k(s,a)}{d^{\pi}(s,a)} = \frac{(1-\gamma)\sum_{t=0}^{\infty}\gamma^t \mathbb{P}(s_{t+k}=s;a_{t+k}=a;\pi)}{(1-\gamma)\sum_{t=0}^{\infty}\gamma^t \mathbb{P}(s_t=s,a_t=a;\pi)}$$
(15)

$$=\frac{\sum_{t=0}^{\infty}\gamma^{t}\mathbb{P}(s_{t+k}=s,a_{t+k}=a;\pi)}{\sum_{t=0}^{\infty}\gamma^{t}\mathbb{P}(s_{t}=s,a_{t}=a;\pi)}$$
(16)

$$\leq \frac{\sum_{t=0}^{\infty} \gamma^t}{\sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s; \pi)} \tag{17}$$

$$= \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s; \pi)$$

$$= \frac{1}{2} \cdot \frac{1}{$$

$$\frac{1}{1-\gamma} \cdot \frac{1}{\sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t=s;\pi)}$$
(18)

where (17) holds by $\mathbb{P}(s_{t+k} = s, a_{t+k} = a; \pi) \leq 1$ and (18) holds by taking the sum of an infinite geometric sequence. By $\sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s, a_t = a; \pi) = \mu_{tar}(s) + \sum_{t=1}^{\infty} \gamma^t \mathbb{P}(s_t = s, a_t = a; \pi)$, we have

$$\frac{1}{1-\gamma} \cdot \frac{1}{\sum_{t=0}^{\infty} \gamma^t \,\mathbb{P}(s_t=s, a_t=a; \pi)} = \frac{1}{1-\gamma} \cdot \frac{1}{\mu(s, a) + \sum_{t=1}^{\infty} \gamma^t \,\mathbb{P}(s_t=s, a_t=a; \pi)}$$
(19)

$$\leq \frac{1}{(1-\gamma)\mu(s,a)} \tag{20}$$

where (20) holds by $\sum_{t=1}^{\infty} \gamma^t \mathbb{P}(s_t = s, a_t = a; \pi) \ge 0.$

B PROOF OF THEOREM 1

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 5. Under Algorithm 1, for any $t \in \mathbb{N}$, we have

$$\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\left(\bar{f}^{t}(s,a) - A^{\pi^{t}}(s,a)\right)^{2}\right] \\
\leq \frac{4}{(1-\gamma)^{2}\mu_{tar,\min}^{2}} \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\left((1-\alpha(t))\epsilon_{td}^{(t)}(s,a) + \alpha(t)\epsilon_{src,be}^{(t)}(s,a;Q_{src})\right)^{2}\right]$$
(21)

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

$$\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\left(\bar{f}^{(t)}(s,a) - A^{\pi^{(t)}}(s,a)\right)^2\right] \tag{22}$$

$$=\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\left(f^{(t)}(s,a) - f^{(t)}(s,\pi^{(t)}(s)) - Q^{\pi^{(t)}}(s,a) + Q^{\pi^{(t)}}(s,\pi^{(t)}(s))\right)^2\right]$$
(23)

$$\leq \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[2 \left(f^{(t)}(s,a) - Q^{\pi^{(t)}}(s,a) \right)^2 + 2 \left(Q^{\pi^{(t)}}(s,\pi^{(t)}(s)) - f^{(t)}(s,\pi^{(t)}(s)) \right)^2 \right]$$
(24)

where (24) holds by the fact that $(x + y)^2 \le 2x^2 + 2y^2$ for any $x, y \in \mathbb{R}$. Then, by linearity of expectation, we obtain

$$\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[2 \left(f^{(t)}(s,a) - Q^{\pi^{(t)}}(s,a) \right)^2 + 2 \left(Q^{\pi^{(t)}}(s,\pi^{(t)}(s)) - f^{(t)}(s,\pi^{(t)}(s)) \right)^2 \right]$$
(25)
$$= \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[2 \left(f^{(t)}(s,a) - Q^{\pi^{(t)}}(s,a) \right)^2 \right] + \mathbb{E}_{s\sim d^{\pi^{(t)}}} \left[2 \left(Q^{\pi^{(t)}}(s,\pi^{(t)}(s)) - f^{(t)}(s,\pi^{(t)}(s)) \right)^2 \right]$$
(26)

$$=\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[2\left(f^{(t)}(s,a)-Q^{\pi^{(t)}}(s,a)\right)^{2}\right]+\mathbb{E}_{s\sim d^{\pi^{(t)}}}\left[2\left(\mathbb{E}_{a'\sim\pi^{(t)}(s)}\left[Q^{\pi^{(t)}}(s,a')-f^{(t)}(s,a')\right]\right]^{2}\right]$$
(27)

$$\leq \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[2 \left(f^{(t)}(s,a) - Q^{\pi^{(t)}}(s,a) \right)^2 \right] + \mathbb{E}_{(s,a')\sim d^{\pi^{(t)}}} \left[2 \left(Q^{\pi^{(t)}}(s,a') - f^{(t)}(s,a') \right)^2 \right]$$
(28)

$$\leq 4\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\left(f^{(t)}(s,a) - Q^{\pi^{(t)}}(s,a)\right)^2\right]$$
(29)

where (28) holds by Jensen's inequality. Then, we proceed to derive an upper bound on $\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\left(f^{(t)}(s,a)-Q^{\pi^{(t)}}(s,a)\right)^2\right]$. By the definition of $f^{(t)} := (1-\alpha(t))Q^{(t)}_{tar}(s,a) + \alpha(t)Q_{src}(\phi^{(t)}(s),\psi^{(t)}(a))$, we have

$$\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\left(f^{(t)}(s,a) - Q^{\pi^{(t)}}(s,a)\right)^2\right]$$
(30)

$$= \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left((1-\alpha(t))Q_{\text{tar}}^{(t)}(s,a) + \alpha(t)Q_{\text{src}}(\phi^{(t)}(s),\psi^{(t)}(a)) - Q^{\pi^{(t)}}(s,a) \right)^2 \right]$$
(31)

$$= \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left(\left(1 - \alpha(t) \right) \left(Q_{\text{tar}}^{(t)}(s,a) - r_{\text{tar}}(s,a) + r_{\text{tar}}(s,a) \right) + r_{\text{tar}}(s,a) \right)^{2} \right]$$

$$+ \alpha(t) \left(Q_{\text{src}}(\phi^{(t)}(s), \psi^{(t)}(a)) - r_{\text{tar}}(s,a) + r_{\text{tar}}(s,a) \right) - Q^{\pi^{(t)}}(s,a) \right)^{2} \right]$$

$$= \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left(\left(1 - \alpha(t) \right) \left(Q_{\text{tar}}^{(t)}(s,a) - r_{\text{tar}}(s,a) + r_{\text{tar}}(s,a) - \gamma \mathbb{E}_{s'\sim P_{\text{tar}}(\cdot|s,a)} [Q_{\text{tar}}^{(t)}(s',a')] \right) + \alpha(t) \left(Q_{\text{src}}(\phi^{(t)}(s), \psi^{(t)}(a)) - r_{\text{tar}}(s,a) + r_{\text{tar}}(s,a) - \gamma \mathbb{E}_{s'\sim P_{\text{tar}}(\cdot|s,a)} [Q_{\text{tar}}^{(t)}(s',a')] \right) + \alpha(t) \left(Q_{\text{src}}(\phi^{(t)}(s), \psi^{(t)}(a)) - r_{\text{tar}}(s,a) + r_{\text{tar}}(s,a) - \gamma \mathbb{E}_{s'\sim P_{\text{tar}}(\cdot|s,a)} [Q_{\text{src}}(\phi^{(t)}(s'), \psi^{(t)}(a'))] + \gamma \mathbb{E}_{s'\sim P_{\text{tar}}(\cdot|s,a)} [Q_{\text{src}}(\phi^{(t)}(s'), \psi^{(t)}(a'))] \right) \\ - Q^{\pi^{(t)}}(s,a) \right)^{2} \right]$$

$$(32)$$

(33)

$$= \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left((1 - \alpha(t)) \left(Q_{\text{tar}}^{(t)}(s,a) - r_{\text{tar}}(s,a) - \gamma \mathbb{E}_{s'\sim P_{\text{tar}}(\cdot|s,a)} [Q_{\text{tar}}^{(t)}(s',a')] \right) \right. \\ \left. + \alpha(t) \left(Q_{\text{src}}(\phi^{(t)}(s), \psi^{(t)}(a)) - r_{\text{tar}}(s,a) - \gamma \mathbb{E}_{s'\sim P_{\text{tar}}(\cdot|s,a)} [Q_{\text{src}}(\phi^{(t)}(s'), \psi^{(t)}(a'))] \right) \right. \\ \left. + (1 - \alpha(t)) \gamma \mathbb{E}_{s'\sim P_{\text{tar}}(\cdot|s,a)} [Q_{\text{tar}}^{(t)}(s',a')] + \alpha(t) \gamma \mathbb{E}_{s'\sim P_{\text{tar}}(\cdot|s,a)} [Q_{\text{src}}(\phi^{(t)}(s'), \psi^{(t)}(a'))] \right. \\ \left. + r_{\text{tar}}(s,a) - Q^{\pi^{(t)}}(s,a) \right)^{2} \right] \\ \left. + r_{\text{tar}}(s,a) - Q^{\pi^{(t)}}(s,a) \right)^{2} \right] \\ \left. = \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left((1 - \alpha(t)) \left(Q_{\text{tar}}^{(t)}(s,a) - r_{\text{tar}}(s,a) - \gamma \mathbb{E}_{s'\sim P_{\text{tar}}(\cdot|s,a)} [Q_{\text{tar}}^{(t)}(s',a')] \right) \right. \\ \left. + \alpha(t) \left(Q_{\text{src}}(\phi^{(t)}(s), \psi^{(t)}(a)) - r_{\text{tar}}(s,a) - \gamma \mathbb{E}_{s'\sim P_{\text{tar}}(\cdot|s,a)} [Q_{\text{src}}^{(t)}(s'), \psi^{(t)}(a'))] \right) \right. \\ \left. + \gamma \mathbb{E}_{s'\sim P_{\text{tar}}(\cdot|s,a)} [f^{(t)}(s',a')] + r_{\text{tar}}(s,a) - Q^{\pi^{(t)}}(s,a) \right)^{2} \right]$$

$$(35)$$

where we obtain (32) by adding the dummy terms $(1 - \alpha(t))(-r_{tar}(s,a) + r_{tar}(s,a))$ and $\alpha(t)(-r_{tar}(s,a) + r_{tar}(s,a))$ to the inner part of (31), (33) 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 $\alpha(t)(-\alpha(t))(-\alpha(t))(s')$ $\gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[Q_{\text{src}}(\phi^{(t)}(s'),\psi^{(t)}(a'))] + \gamma \mathbb{E}_{s' \sim P_{\text{tar}}(\cdot|s,a)}[Q_{\text{src}}(\phi^{(t)}(s'),\psi^{(t)}(a'))])$ to the inner part $\alpha' \sim \alpha^{(t)}(\cdot|s')$ $a' \sim \pi^{(t)}(\cdot|s')$ $a' \sim \pi^{(t)}(\cdot|s')$ of (32), (34) holds by rearranging the terms in (33), and (35) holds by the definition of $f^{(t)}$. Then, by adding $\gamma \mathbb{E}_{s'' \sim P_{tar}(\cdot|s,a)} [Q^{\pi^{(t)}}(s'',a'')] - \gamma \mathbb{E}_{s'' \sim P_{tar}(\cdot|s,a)} [Q^{\pi^{(t)}}(s'',a'')]$ to the inner part of (35), we $a'' \sim \pi^{(t)}(\cdot|s'')$ can rewrite (35) as

$$\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left((1-\alpha(t)) \left(Q_{tar}^{(t)}(s,a) - r_{tar}(s,a) - \gamma \mathbb{E}_{s'\sim P_{tar}(\cdot|s,a)} [Q_{tar}^{(t)}(s',a')] \right) \right. \\ \left. + \alpha(t) \left(Q_{src}(\phi^{(t)}(s),\psi^{(t)}(a)) - r_{tar}(s,a) - \gamma \mathbb{E}_{s'\sim P_{tar}(\cdot|s,a)} [Q_{src}(\phi^{(t)}(s'),\psi^{(t)}(a'))] \right) \right. \\ \left. + \gamma \mathbb{E}_{s'\sim P_{tar}(\cdot|s,a)} [f^{(t)}(s',a')] + r_{tar}(s,a) - Q^{\pi^{(t)}}(s,a) \right. \\ \left. + \gamma \mathbb{E}_{s'\sim \pi^{(t)}(\cdot|s')} \left[Q^{\pi^{(t)}}(s'',a'') \right] - \gamma \mathbb{E}_{s''\sim P_{tar}(\cdot|s,a)} [Q^{\pi^{(t)}}(s'',a'')] \right)^2 \right] \\ \left. + \gamma \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left| (1-\alpha(t)) \left(Q_{tar}^{(t)}(s,a) - r_{tar}(s,a) - \gamma \mathbb{E}_{s'\sim P_{tar}(\cdot|s,a)} [Q_{tar}^{(t)}(s',a')] \right) \right. \\ \left. + \alpha(t) \left(Q_{src}(\phi^{(t)}(s),\psi^{(t)}(a)) - r_{tar}(s,a) - \gamma \mathbb{E}_{s'\sim P_{tar}(\cdot|s,a)} [Q_{tar}(s',a')] \right) \right. \\ \left. + \gamma \mathbb{E}_{s'\sim P_{tar}(\cdot|s,a)} [f^{(t)}(s',a')] + r_{tar}(s,a) - Q^{\pi^{(t)}}(s,a) \right] \right] \right. \\ \left. + \gamma \mathbb{E}_{s'\sim P_{tar}(\cdot|s,a)} [f^{(t)}(s',a')] + r_{tar}(s,a) - Q^{\pi^{(t)}}(s,a) \right] \right]$$

$$\left. + \gamma \mathbb{E}_{s'\sim P_{tar}(\cdot|s,a)} [Q^{\pi^{(t)}}(s'',a'')] - \gamma \mathbb{E}_{s''\sim P_{tar}(\cdot|s,a)} [Q^{\pi^{(t)}}(s'',a'')] \right]^2 \right]$$

$$\left. + \gamma \mathbb{E}_{s''\sim P_{tar}(\cdot|s,a)} [Q^{\pi^{(t)}}(s'',a'')] - \gamma \mathbb{E}_{s''\sim P_{tar}(\cdot|s,a)} [Q^{\pi^{(t)}}(s'',a'')] \right]^2 \right]$$

$$\left. + \gamma \mathbb{E}_{s''\sim P_{tar}(\cdot|s,a)} [Q^{\pi^{(t)}}(s'',a'')] - \gamma \mathbb{E}_{s''\sim P_{tar}(\cdot|s,a)} [Q^{\pi^{(t)}}(s'',a'')] \right]^2 \right]$$

$$\leq \mathbb{E}_{(s,a)\sim d^{\pi^{(1)}}} \left[\left(\left| (1 - \alpha(t)) \left(Q_{tar}^{(tr)}(s,a) - r_{tar}(s,a) - \gamma \mathbb{E}_{s'\sim P_{tar}(\cdot|s,a)}^{(1)} \left[Q_{tar}^{(t)}(s',a') \right] \right) \right| \\ + \left| \alpha(t) \left(Q_{scc}(\phi^{(t)}(s), \psi^{(t)}(a)) - r_{tar}(s,a) - \gamma \mathbb{E}_{s'\sim P_{tar}(\cdot|s,a)}^{(1)} \left[Q_{tar}(\phi^{(t)}(s'), \psi^{(t)}(a')) \right] \right) \right| \\ + \left| \gamma \mathbb{E}_{s'\sim P_{tar}(\cdot|s,a)}^{(1)} \left[f^{(t)}(s',a') \right] + r_{tar}(s,a) - Q^{\pi^{(t)}}(s,a) \right| \\ + \left| \gamma \mathbb{E}_{s'\sim \pi^{(t)}(\cdot|s')}^{(1)} \left[Q^{(t)}(s',a'') \right] - \gamma \mathbb{E}_{s'\sim P_{tar}(\cdot|s,a)}^{(1)} \left[Q^{(t)}(s'',a'') \right] \right] \right)^{2} \right] \\ \leq \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left((1 - \alpha(t)) \left| Q_{tar}^{(t)}(s,a) - r_{tar}(s,a) - \gamma \mathbb{E}_{s'\sim P_{tar}(\cdot|s,a)} \left[Q_{tar}^{(t)}(s',a') \right] \right| \right] \\ + \alpha(t) \left| \left(Q_{src}(\phi^{(t)}(s), \psi^{(t)}(a)) - r_{tar}(s,a) - \gamma \mathbb{E}_{s'\sim P_{tar}(\cdot|s,a)} \left[Q_{tar}^{(t)}(s',a') \right] \right| \right] \\ = \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}(\cdot|s')} \left[f^{(t)}(s',a') - \gamma \mathbb{E}_{s'\sim P_{tar}(\cdot|s,a)} \left[Q_{src}(\phi^{(t)}(s'), \psi^{(t)}(a')) \right] \right] \right] \\ + \left| \gamma \mathbb{E}_{s'\sim P_{tar}(\cdot|s,a)} \left[f^{(t)}(s',a') - \gamma \mathbb{E}_{s'\sim P_{tar}(\cdot|s,a)} \left[Q_{src}(\phi^{(t)}(s'), \psi^{(t)}(a')) \right] \right] \right] \\ + \left| \gamma \mathbb{E}_{s'\sim P_{tar}(\cdot|s,a)} \left[f^{(t)}(s',a') - \gamma \mathbb{E}_{s'\sim P_{tar}(\cdot|s,a)} \left[Q_{src}(\phi^{(t)}(s'), \psi^{(t)}(a')) \right] \right] \right]^{2} \right] \\ = \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left((1 - \alpha(t)) e_{td}^{(t)}(s,a) + \alpha(t) e_{src,hc}^{(t)}(s',a'') \right] \right] \right]^{2} \\ = \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[f^{(t)}(s',a') - Q^{\pi^{(t)}}(s',a') \right] \right] \right]^{2} \\ = \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left[f^{(t)}(s',a') - Q^{\pi^{(t)}}(s',a') \right] \right] \right]^{2} \\ = \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left[f^{(t)}(s',a') - Q^{\pi^{(t)}}(s',a') \right] \right] \right]^{2} \\ = \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left[f^{(t)}(s',a') - Q^{\pi^{(t)}}(s',a') \right] \right]^{2} \\ = \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left[f^{(t)}(s',a') - Q^{\pi^{(t)}}(s',a') \right] \right]^{2} \\ = \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left| f^{(t)}(s',a') - Q^{\pi^{(t)}}(s',a') \right| \right]^{2} \\ = \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left| f^{(t)}(s',a') - Q^{\pi^{(t)}}(s',a') \right| \right]^{2} \\ = \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left| f^{(t)}(s',a') - Q^{\pi^{(t)}}(s',a') \right| \right]^{2} \\ = \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left| f^{(t)}(s',a') - Q^{\pi^{(t)}}(s',a')$$

where (37) holds by the fact that $x^2 = |x|^2$, (38) holds by triangle inequality, (39) by the facts that $0 \le \alpha(t) \le 1$ and $0 \le 1 - \alpha(t) \le 1$, (40) holds by coupling (s', a') and (s'', a'') and applying Bellman expectation equation as well as the definitions that $\epsilon_{td}^{(t)}(s,a) := |Q_{tar}^{(t)}(s,a) - r_{tar}(s,a) - \gamma \mathbb{E}_{s' \sim P_{tar}(\cdot|s,a)}[Q_{tar}^{(t)}(s',a')]|$ and $\epsilon_{src,be}^{(t)}(s,a;Q_{src}) := |Q_{src}(\phi^{(t)}(s),\psi^{(t)}(a)) - r_{tar}(s,a) - a' \sim \pi^{(t)}(\cdot|s')$ $\gamma \mathbb{E}_{s' \sim P_{tar}(\cdot|s,a)} [Q_{src}(\phi^{(t)}(s'), \psi^{(t)}(a'))]|$. By recursively applying the procedure from (30) to (41) $a' \sim \pi^{(t)}(\cdot | s')$ to $|f^{(t)}(s',a') - Q^{\pi^{(t)}}(s',a')|$, 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:

$$\mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}}\left[\left(f^{(t)}(s,a) - Q^{\pi^{(t)}}(s,a)\right)^2\right]$$
(42)

(43)

$$\leq \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left(\left(1-\alpha(t)\right)\epsilon_{\mathrm{td}}^{(t)}(s,a) + \alpha(t)\epsilon_{\mathrm{src,be}}^{(t)}(s,a;Q_{\mathrm{src}}) \right) \right]^{2} \right]$$

970
971
$$+ \gamma \mathbb{E}_{s' \sim P_{tar}(\cdot | s, a)} \left[\left| f^{(t)}(s', a') - Q^{\pi^{(t)}}(s', a') \right| \right] \right)^2 \right]$$

$$a' \sim \pi^{(t)}(\cdot | s')$$

$$\leq \mathbb{E}_{(s,a)\sim d^{\pi(t)}} \left[\left(\left(1 - \alpha(t)\right) \epsilon_{td}^{(t)}(s,a) + \alpha(t) \epsilon_{src,be}^{(t)}(s,a;Q_{src}) + \frac{1}{(1 - \gamma)\mu_{tar,min}} \left(\gamma \left(1 - \alpha(t)\right) \epsilon_{td}^{(t)}(s,a) + \gamma \alpha(t) \epsilon_{src,be}^{(t)}(s,a;Q_{src}) \right) \right] \right]$$

$$(45)$$

(44)

$$+ \gamma^2 (1 - \alpha(t)) \epsilon_{td}^{(t)}(s, a) + \gamma^2 \alpha(t) \epsilon_{src, be}^{(t)}(s, a; Q_{src}) + \cdots) \bigg)^2 \bigg]$$

 $\leq \mathbb{E}_{(s,a)\sim d^{\pi(t)}} \left| \left(\left(1-\alpha(t)\right) \epsilon_{\mathrm{td}}^{(t)}(s,a) + \alpha(t) \epsilon_{\mathrm{src,be}}^{(t)}(s,a;Q_{\mathrm{src}}) \right) \right.$

 $+ \gamma \mathbb{E}_{\substack{s' \sim P_{\text{tar}}(\cdot|s,a) \\ a' \sim \pi^{(t)}(\cdot|s')}} \Big[(1 - \alpha(t)) \epsilon_{\text{td}}^{(t)}(s',a') + \alpha(t) \epsilon_{\text{src,be}}^{(t)}(s',a';Q_{\text{src}}) \Big]$

 $+ \mathbb{E}_{\substack{s'' \sim P_{\text{tar}}(\cdot | s', a') \\ a'' \sim \pi^{(t)}(\cdot | s'')}} \Big[\Big| f^{(t)}(s'', a'') - Q^{\pi^{(t)}}(s'', a'') \Big| \Big] \Big] \Big)^2 \Big]$

$$\leq \frac{1}{(1-\gamma)^4 \mu_{\text{tar,min}}^2} \mathbb{E}_{(s,a)\sim d^{\pi^{(t)}}} \left[\left((1-\alpha(t))\epsilon_{\text{td}}^{(t)}(s,a) + \alpha(t)\epsilon_{\text{src,be}}^{(t)}(s,a;Q_{\text{src}}) \right)^2 \right]$$
(46)

where (44) holds by applying the procedure from (30) to (41) to $f^{(t)}(s', a') - Q^{\pi^{(t)}}(s', a')$, (45) holds by applying the procedure from (30) to (41) to all the subsequent time steps and using importance sampling with the importance ratio bound in Lemma 4 and then using the same dummy variables (s, a)for all the subsequent state-action pairs, and (46) holds by taking the sum of an infinite geometric sequence. \square

Theorem 1. (Average Sub-Optimality) Under the QAvatar in Algorithm 1 and Assumption 1, given any fixed learning rate $\eta > 0$, the average sub-optimality over T iterations can be upper bounded as

$$\frac{1}{T} \sum_{t=1}^{T} \left(V^{\pi^*}(\mu_{tar}) - V^{\pi^{(t)}}(\mu_{tar}) \right) \leq \underbrace{\frac{1}{(1-\gamma)T} \sum_{t=1}^{T} \mathbb{E}_{(s,a)\sim d^{\pi^*}} \left[\max_{a'} \bar{f}^{(t)}(s,a') \right] + \frac{\log |\mathcal{A}_{tar}|}{(1-\gamma)T\eta}}_{(s,a')}$$

$$+\underbrace{\frac{C_{0}}{T}\sum_{t=1}^{T}\alpha(t)\|\epsilon_{src, be}^{(t)}(Q_{src})\|_{\infty}}_{(b)}+\underbrace{\frac{C_{0}}{T}\sum_{t=1}^{T}(1-\alpha(t))\|\epsilon_{td}^{(t)}\|_{\infty}}_{(c)},$$
(4)

where
$$C_0 := 2\sqrt{C_{\pi^*}}/((1-\gamma)^2\mu_{tar, min})$$
 and $\bar{f}^{(t)}(s, a) := f^{(t)}(s, a) - f^{(t)}(s, \pi^{(t)}(s)).$

Proof. We start by providing an upper bound on the sub-optimality gap $V^{\pi^*}(\mu_{tar}) - V^{\pi^{(t)}}(\mu_{tar})$ at each iteration. Recall that d_{tar}^{π} denotes the discounted state-action visitation distribution of policy π in the target domain. Note that

$$V^{\pi^{*}}(\mu_{\rm tar}) - V^{\pi^{(t)}}(\mu_{\rm tar})$$
(47)

$$= \frac{1}{1 - \gamma} \mathbb{E}_{(s,a) \sim d_{\text{tar}}^{\pi^*}} \left[A^{\pi^{(t)}}(s,a) \right]$$
(48)

$$= \frac{1}{1-\gamma} \mathbb{E}_{(s,a)\sim d_{\text{tar}}^{\pi^*}} \left[\bar{f}^{(t)}(s,a) - \bar{f}^{(t)}(s,a) + A^{\pi^{(t)}}(s,a) \right]$$
(49)

$$= \frac{1}{1 - \gamma} \mathbb{E}_{(s,a) \sim d_{\text{tar}}^{\pi^*}} \left[\bar{f}^{(t)}(s,a) \right] + \frac{1}{1 - \gamma} \mathbb{E}_{(s,a) \sim d_{\text{tar}}^{\pi^*}} \left[-\bar{f}^{(t)}(s,a) + A^{\pi^{(t)}}(s,a) \right]$$
(50)

$$\leq \frac{1}{1-\gamma} \mathbb{E}_{(s,a)\sim d_{\text{tar}}^{\pi^*}} \left[\bar{f}^{(t)}(s,a) \right] + \frac{1}{1-\gamma} \sqrt{\mathbb{E}_{(s,a)\sim d_{\text{tar}}^{\pi^*}} \left[\left(-\bar{f}^{(t)}(s,a) + A^{\pi^{(t)}}(s,a) \right)^2 \right]},$$
 (51)

where (48) holds by the performance difference lemma (cf. Lemma 1), (49) is obtained by adding $f^t(s, a) - f^t(s, a)$, (50) is obtained by rearranging the terms in (49), and (51) holds by Jensen's

inequality. By the fact that $\|\frac{d^{\pi^*}}{d\pi^{(t)}}\|_{\infty} \leq C$, we have 1027 1028 $\frac{1}{1-\gamma} \mathbb{E}_{(s,a)\sim d^{\pi^*}} \left[\bar{f}^{(t)}(s,a) \right] + \frac{1}{1-\gamma} \sqrt{\mathbb{E}_{s,a\sim d^{\pi^*}} \left[\left(-\bar{f}^{(t)}(s,a) + A^{\pi^{(t)}}(s,a) \right)^2 \right]}$ (52)1029 1030 $\leq \frac{1}{1-\gamma} \mathbb{E}_{(s,a)\sim d^{\pi^*}} \left[\bar{f}^{(t)}(s,a) \right] + \frac{1}{1-\gamma} \sqrt{C \cdot \mathbb{E}_{s,a\sim d^{\pi^{(t)}}} \left[\left(-\bar{f}^{(t)}(s,a) + A^{\pi^{(t)}}(s,a) \right)^2 \right]}.$ 1031 (53) 1032 1033 Recall the definitions of $\epsilon_{td}^{(t)}(s,a)$ and $\epsilon_{src}^{(t)}{}_{be}(s,a;Q_{src})$ as 1034 $\epsilon_{\mathrm{td}}^{(t)}(s,a) := \big| Q_{\mathrm{tar}}^{(t)}(s,a) - r_{\mathrm{tar}}(s,a) - \gamma \mathbb{E}_{s' \sim P_{\mathrm{tar}}(\cdot|s,a)}[Q_{\mathrm{tar}}^{(t)}(s',a')] \big|,$ $a' \sim \pi^{(t)}(\cdot|s')$ 1035 (54)1036 1037 $\epsilon_{\rm src,be}^{(t)}(s,a;Q_{\rm src}) := \big| Q_{\rm src}(\phi^{(t)}(s),\psi^{(t)}(a)) - r_{\rm tar}(s,a) - \gamma \mathbb{E}_{s' \sim P_{\rm tar}(\cdot|s,a)}[Q_{\rm src}(\phi^{(t)}(s'),\psi^{(t)}(a'))] \big|.$ 1038 1039 (55)1040 1041 Recall that we also define 1042 $\|\epsilon_{\mathsf{td}}^{(t)}\|_{\infty} := \max_{(s,a)\in\mathcal{S}\times\mathcal{A}} \epsilon_{\mathsf{td}}^{(t)}(s,a),$ (56)1043 $\|\epsilon_{\mathrm{src, be}}^{(t)}(Q_{\mathrm{src}})\|_{\infty} := \max_{(s,a)\in\mathcal{S}\times\mathcal{A}} \epsilon_{\mathrm{src, be}}^{(t)}(s,a;Q_{\mathrm{src}}).$ 1044 (57)1045 1046

1047 We are ready to put everything together and establish the cumulative sub-optimality. By taking the 1048 summation of (53) over iterations, we have 1049

$$\sum_{t=1}^{T} \mathbb{E}_{s \sim \mu_{\text{tar}}} \Big[V^{\pi^*}(s) - V^{\pi^{(t)}}(s) \Big]$$
(58)

$$\leq \sum_{t=1}^{T} \frac{1}{1-\gamma} \mathbb{E}_{(s,a)\sim d^{\pi^*}} \left[\bar{f}^{(t)}(s,a) \right] + \sum_{t=1}^{T} \frac{1}{1-\gamma} \sqrt{C \mathbb{E}_{s,a\sim d^{\pi^{(t)}}} \left[\left(-\bar{f}^{(t)}(s,a) + A^{\pi^{(t)}}(s,a) \right)^2 \right]}$$
(59)

1058

1061 1062

1063 1064

1067

1068

1026

$$\leq \sum_{t=1}^{T} \frac{1}{1-\gamma} \mathbb{E}_{(s,a)\sim d^{\pi^{*}}} \left[\max_{a'} \bar{f}^{(t)}(s,a') \right] + \frac{\log |\mathcal{A}_{tar}|}{(1-\gamma)\eta} + \frac{2\sqrt{C}}{(1-\gamma)^{2}\mu_{tar,\min}} \sum_{t=1}^{T} \alpha(t) \|\epsilon_{src, be}^{(t)}(Q_{src})\|_{\infty} + \frac{2\sqrt{C}}{(1-\gamma)^{2}\mu_{tar,\min}} \sum_{t=1}^{T} (1-\alpha(t)) \|\epsilon_{td}^{(t)}\|_{\infty}.$$
(60)

where (59) follows directly from (53) and (60) holds by Lemma 5 and Lemma 2.

С Additional Experimental Results

C.1 TOY EXAMPLE TO SHOW THE BENEFIT OF CROSS-DOMAIN BELLMAN-LIKE LOSS.

We consider the 3-by-3 grid navigation problem, 1069 as shown in Figure 5. In both domains, there are 1070 only two actions: 'going top' and 'going right.' 1071 The state of the source domain is described in 1072 decimal coordinates, while the state of the target domain is described in binary coordinates. The 1074 white squares represent obstacles that cannot 1075 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 1077 1078 at this state, and the agent will receive an ending reward of +1. (iii) Treasure state: When the 1079 agent first navigates to this state, it will receive



Figure 5: The source and target domain of the grid navigation example.

1080 +0.5 rewards. In other states or at other times navigating the treasure state, the agent will not receive 1081 any reward. In the source domain, the start state, end state, and treasure state are set to (0,0), (0,2), 1082 and (2,2), respectively. In the target domain, the start state, end state, and treasure state are set to 1083 (0, 0, 0, 0), (0, 0, 1, 1),and (1, 1, 1, 1), respectively. We assume that the source Q-function Q_{src} is 1084 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 $(0,0) \rightarrow (0,1) \rightarrow (0,2) \rightarrow (1,2) \rightarrow (2,2)$ and the optimal trajectory of the target domain is $(0,0,0,0) \rightarrow (0,0,0,1) \rightarrow (0,0,1,1) \rightarrow (0,1,1,1) \rightarrow (0,1,1) \rightarrow (0,1,1,1) \rightarrow (0,1,1) \rightarrow (0,1) \rightarrow ($ 1086 (1, 1, 1, 1). Consider two trajectories in the source domain: Traj-A, which is the optimal trajectory, 1087 and Traj-B, defined as $(0,0) \rightarrow (0,1) \rightarrow (1,1) \rightarrow (1,2) \rightarrow (2,2)$. When we map the optimal 1088 trajectory of the target domain to Traj-A and the optimal trajectory of the target domain to Traj-B, 1089 both mappings result in 0 cycle consistency loss. This suggests that the cycle consistency cannot 1090 determine which mapping is superior. This phenomenon results from the unsupervised nature of 1091 dynamics cycle consistency. In contrast, when we mapping the optimal trajectory of the target domain 1092 to Traj-A yields a cross-domain Bellman-like loss of 0, while mapping the optimal trajectory of 1093 the target domain to Traj-B results in a cross-domain Bellman-like loss of 1. Thus, we can achieve 1094 optimal mapping results based on the cross-domain Bellman error, while the cycle consistency loss 1095 provides sub-optimal mapping results.

1096 1097

1098

1109

1125 1126

1127

C.2 REPRODUCTION AND SANITY CHECKS FOR DCC, CMD, AND CAT

1099 In this section, we report the reproduced scores of DCC, CMD, and CAT as sanity checks for success-1100 ful reproduction. To simplify the expressions, we abbreviate 'original' as 'orig' and 'reproduced' as 1101 'repr'. As shown in Table 2, we can observe that QAvatar is indeed outperform than DCC, CMD. For 1102 the CAT, since the original paper only provides the result of 2-to-1 transfer (i.e., two source domains 1103 and one target domain) from CentipedeFour and CentipedeEight to CentipedeSix and get around 2200 episodic return at timesteps 700000. Also, in the original paper, CAT has reproduced the MIKT, 1104 which is equivalent to CAT for 1-to-1 transfer (i.e., one source domain and one target domain), and 1105 get around 1900 episodic rewards at timesteps 700000. In our reproduction, CAT achieves 1715 1106 evaluation return, and this is close to MIKT in the original CAT paper. Hence, we believe that our 1107 implementation of CAT is correct given that our configuration is the same as MIKT. 1108

Table 2: The original and reproduced results of DCC and CMD compared with QAvatar and SAC.

Environment	DCC (orig)	DCC (repr)	CMD (orig)	CMD (repr)	SAC	QAvatar
Swimmer Halfcheetah Ant	$\begin{array}{c} 204\pm 56\\ 2471\pm 382\\ \text{N/A} \end{array}$	$\begin{array}{c} 132 \pm 47 \\ 1360 \pm 729 \\ 973 \pm 501 \end{array}$	$\begin{array}{c} 44 \pm 38 \\ 2114 \pm 332 \\ 649 \pm 347 \end{array}$	$\begin{array}{c} 41 \pm 14 \\ 303 \pm 75 \\ 882 \pm 52 \end{array}$	$\begin{array}{c} 235 \pm 141 \\ 11445 \pm 1897 \\ 2290 \pm 785 \end{array}$	$\begin{array}{c} 316\pm139\\ 12819\pm679\\ 2840\pm1532 \end{array}$

1117 C.3 FINAL REWARDS

In this section, we show the asymptotic performance of all baselines and our algorithm. In the MuJoCo environments except for Ant and Inverted Double Pendulum, we train all the target-domain models for 500k steps. In Ant and Inverted Double Pendulum, we train all the target-domain models for 350k and 20k steps, respectively. In Robosuite environments, we train all the target-domain models for 20k steps. The asymptotic performances of all baselines and our algorithm are shown in the following tables.

Table 3: Final rewards of QAvatar and all baselines in the MuJoCo environments.

28	Algorithm	Hopper	HalfCheetah	Ant	Centipede	Reacher	Modified IDP
9	QAvatar	$\textbf{2762} \pm \textbf{440}$	12316 ± 586	$\textbf{2234} \pm \textbf{1112}$	$\textbf{2020} \pm \textbf{1465}$	-6.1 ± 0.3	9241 ± 62
30	SAC	2086 ± 257	10986 ± 1822	1620 ± 527	872 ± 36	$\textbf{-5.5}\pm\textbf{0.1}$	9212 ± 152
1	CMD	59 ± 46	-253 ± 344	778 ± 144	834 ± 6116	-14.8 ± 0.5	72 ± 12
	DCC	30 ± 16	-631 ± 185	-1240 ± 838	148 ± 182	-16.5 ± 1.4	95 ± 6
-	CAT	154 ± 156	46 ± 250	17 ± 27	1715 ± 430	$\textbf{-12.2}\pm1.0$	41 ± 10
00	FT	2530 ± 456	12016 ± 1052	1740 ± 642	1123 ± 508	$\textbf{-5.8}\pm0.2$	$\textbf{9349} \pm \textbf{86}$

1150					
1137	Environment	Block Lifting	Door Opening	Table Wiping	Highway
1138	QAvatar	$\textbf{98.0} \pm \textbf{21.1}$	$\textbf{185.2} \pm \textbf{66.9}$	$\textbf{67.1} \pm \textbf{9.1}$	$132.2{\pm}~31.7$
1139	SAC	90.3 ± 23.4	160.1 ± 40.3	47.2 ± 7.1	117.8 ± 15.0
1140	CMD	0.9 ± 0.6	7.8 ± 6.4	0.8 ± 0.4	13.0 ± 4.7
	DCC	0.6 ± 0.2	8.2 ± 4.7	0.9 ± 0.7	18.1 ± 19.4
1141	CAT	15.0 ± 14.3	34.7 ± 8.4	55.5 ± 29.7	70.1 ± 6.3
1142	FT	21.9 ± 7.8	129.2 ± 44.9	36.8 ± 17.2	119.1 ± 21.1

Table 4: Final performances of *Q*Avatar and all baselines in the Robosuite and BARK-ML environments.

1145 C.4 TIME TO THRESHOLD

In the following table, we discover that QAvatar uses the less data to reach the threshold than SAC does. In Hopper, QAvatar only needs half the amount of data SAC needs to reach the goal.

Table 5: Time to threshold of *Q*Avatar and all baselines

1151					
1152	Environment	Threshold	QAvatar	SAC	SAC / QAvatar
1153	Hopper	2300	252K	836K	3.32
1154	HalfCheetah	10000	288K	400K	1.39
	Ant	1600	254K	344K	1.35
1155	Centipede	900	210K	988K	4.70
1156	Block Lifting	85	90K	94K	1.04
1157	Door Opening	150	80K	94K	1.18
1150	Table Wiping	45	74K	96K	1.30
1158	Highway	110	236K	374K	1.58
1159	Reacher	-7	150K	84K	0.56
1160	Inverted Double Pendulum	9000	16K	18K	1.13

1161 1162

1176

1181

1183

1186

1143 1144

1146

1147

1148 1149

1150

1163 C.5 Ablation Study: Deactivating the Flow Model

1164 As mentioned above, we use a normalizing flow 1165 model to restrict the output range of the mapping functions in the feasible regions. In this 1166 experiment, we disable the flow model and eval-1167 uate QAvatar in Swimmer and Door Opening. 1168 In Figure 6, QAvatar without a flow model per-1169 forms worse than QAvatar with a flow model. 1170 In Door Opening, although the ewma values of 1171 rewards obtained by QAvatar without the flow 1172 model are higher than 100, it has to spend more 1173 time attaining high rewards than QAvatar with 1174 the flow model does. 1175



QAvata	r w/o flow model	🔶 QAvatar	🕂 SAC

Figure 6: Ablation Study: *Q*Avatar without the flow model

1177 D IMPLEMENTATION

1178 DETAILS OF QAVATAR

1179		
1180	D.1	Pseudo Code of the Practical Implementation of $QAvatar$

- 1182 In this section, we provide the pseudo code of the practical version of QAvatar in Algorithm 2.
- 1184 D.2 INTER-DOMAIN MAPPING NETWORK AUGMENTED WITH A NORMALIZING FLOW MODEL
- 1187 As mentioned in Section 4, we use the flow model to map the outputs of the mapping functions to the feasible regions. The way to integrate these two components is shown in Figure 7.



As mentioned in Section 5, the source domains of our experiments are the original MuJoCo environments such as Swimmer-v3, Hopper-v3, HalfCheetah-v3 and Ant-v3. The target domains are the modified MuJoCo environments such as Swimmer with four limbs, Hopper with an extra thigh, HalfCheetah with three legs and Ant with five legs. For the Centipede, CentipedeFour refers to a Centipede with four legs, and CentipedeSix refers to a Centipede with six legs. The environments are shown in Figure 8.

1242

1258

1259

1260

1261 1262 1263

1264

1288

1289

1290 1291

1293

1243 1244 1245 1246 1247 1248 1249 (a) Swimmer (b) Hopper (c) HalfCheetah (d) Ant (e) CentipedeFour 1250 1251 1252 1253 1254 1255 1256 (f) Four-limb Swimmer 1257 (g) Three-thigh Hopper (h) Three-leg HalfCheetah (i) Five-leg Ant (j) CentipedeSix

Figure 8: The environments of the source domains and the target domains. (a)-(e): Source domains - Original MuJoCo environments and CentipedeFour. (f)-(j): Target domains - Modified MuJoCo environments and CentipedeSix.

E.2 **ROBOSUITE AND BARK-ML ENVIRONMENTS**

1265 Robosuite is a popular robot learning package. We evaluate QAvatar on three tasks, including block 1266 lifting, door opening, and table wiping. For each task, we consider cross-domain transfer from 1267 controlling a Panda robot arm to controlling a UR5e robot arm. For the BARK-ML environments, we consider transfer from "merging-v0" to "highway-v0". These four tasks are illustrated in Figure 9. 1268 1269



Figure 9: The environments of the source domains and the target domains. (a)-(d)The source domains: control Panda to solve the tasks in robosuite and merging-v0 in bark-ml. (e)-(h)The target domains: control UR5e to solve the tasks in robosuite and highway-v0 in bark-ml.

E.3 THE IMPLEMENTATION DETAILS OF BASELINES

SAC. The implementation of SAC used in our experiments is released by stable-baselines3 Raffin 1294 et al. (2021). The settings of all hyperparameters except for the discouted factor γ follows the default 1295 settings of SAC in the documentation of stable-baselines3. The discouted factor is set as 0.9999

in Swimmer-v3 and 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. We implement CMD by ourselves according to the pseudocode of CMD shown 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.

DCC. We use the original implementation of (Zhang et al., 2021) (https://github.com/ sjtuzq/Cycle_Dynamics) with their default setting Zhang et al. (2021). For a fair comparison, we use the same source model used in QAvatar and change DCC from collecting the fixed amount of data to collecting data continuously.

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.

Table 6: A list of candidate hyperparameters for Robosuite and MuJoCo.

1327			
1328	Parameter	MuJoCo	Robosuite
1329	learning rate	0.0001, 0.0003, 0.0004, 0.0008	0.0001, 0.0003
1330	length of rollouts	500, 2000 (50, 100 for Modified IDP)	2000
1331	batch size	50, 100 (20, 25 for Modified IDP)	50, 100, 200
1332	entropy coefficient (ent. coef.)	0.01, 0.002	0.01, 0.002
1333	n epochs	10, 20	5, 10
1334	num. of hidden layer of encoder/decoder	1	1
1335	num. of hidden layer of actor/critic	2	2
1336	hidden layer size	256	256

Table	7.	Final	hyper	narameters	chosen	for	each	enviror	ment
rabic	1.	1 mai	nyper	parameters	chosen	101	caci	CIIVIIOI	micin

	learning rate	len. of rollouts	batch size	ent. coef.	n epochs
Hopper	0.0008	2000	100	0.002	20
HalfCheetah	0.0001	500	50	0.002	10
Ant	0.0004	500	50	0.002	10
CentipedeSix	0.0003	2048	64	0.00	10
InvertedDoublePendulum	0.001	100	20	0.01	20
Reacher	0.0003	2048	64	0.00	10
Robosuite	0.0003	2000	100	0.01	10
Highway	0.0003	2048	64	0.00	10

E.4 DETAILED CONFIGURATION OF QAVATAR

The base algorithm, SAC, is implemented by stable-baselines3 Raffin et al. (2021). As for the compute resource, we use NVIDIA GeForce RTX 3090 to do the experiments. Finishing the whole training process including training the source-domain model, target-domain model and flow model once needs 44 hours in the MuJoCo environments and 39 hours in the Robosuite environments. The Hyperparameters of QAvatar are shown in the following two tables. The consider functions of decay functions are $1/\sqrt{t}$, 1/t, $1/t^2$ and $1/t^3$ and the final decay functions chosen for each environments are shown in the table 9. The settings of hyperparameters such as critic/actor learning rate, batch size, buffer size and discounted factor are same as SAC.

1362		
1363		
1364		
1365		
1366		
1367		
1368	Table 8: A list of hyperparameters of QAv	atar .
1369		
1370	Parameter	Value
1371	critic/actor learning rate	0.0003
1372	state mapping function learning rate	0.0005
1373	action mapping function learning rate	0.01
1374	batch size	256
1375	replay buffer size	10^{6}
1376	optimizer	Adam
1377	number of hidden layer of mapping functions	1
1378	hidden layer size	256
1379		
1380		
1381		
1382		
1383		
1384		
1385		
1386		
1387		
1388		
1389		
1390		
1391		604
1392	Table 9: A list of environment-specific hyperparamet	ers of QAV
1393		
1394	Environment Decay Func	α
1395	Hopper-v3 $1/t^2$	
1396	HalfCheetah-v3 $1/t^3$	
1397	Ant-v3 $1/\sqrt{t}$	
1398	CentipedeSix $1/t$	
1399	InvertedDoublePendulum-v2 $1/t$	
1400	Reacher-v2 $1/t^3$	
1401	Block Lifting $1/t$	
1402	Door Opening $1/t$	
1403	Table wiping $1/t$	
1-100	Highway-v0 $1/t^3$	

vatar .







Figure F.4: Training curves of QAvatar and
SAC in the negative transfer scenario of locomotion tasks.

Figure F.5: Training curves of QAvatar and SAC in the opposite Q-function of source and target domain Transfer Scenario.



Figure F.3: Training Curves of SAC and *Q*Avatar in three different setting, where "src1-tar" refer to transfer domain src1 to domain tar.(src1: Ant-v3 with front left and back right legs disabled, src2: Ant-v3 with front right and back left Legs disable, tar: original Ant-v3).





Figure F.7: Training Curves of *Q*Avatar, SAC, and PAR in Various Environments Under the Same Settings as in Section 5.