

IOB: Integrating Optimization Transfer and Behavior Transfer for Multi-Policy Reuse

Siyuan Li¹, Hao Li², Jin Zhang³, Zhen Wang², Peng Liu¹,
Chongjie Zhang^{3*}

¹Faculty of Computing, Harbin Institute of Technology, Harbin, 150001, China.

²School of Artificial Intelligence, Optics and Electronics (iOPEN), Northwestern Polytechnical University, Xi'an, 710072, China.

^{3*}McKelvey School of Engineering, Washington University in St. Louis, St. Louis, 63130, United States.

*Corresponding author(s). E-mail(s): chongjie@wustl.edu;

Contributing authors: siyuanli@hit.edu.cn; li.hao@mail.nwpu.edu.cn;
jin-zhan20@mails.tsinghua.edu.cn; zhenwang0@gmail.com;
pengliu@hit.edu.cn;

Abstract

Humans have the ability to reuse previously learned policies to solve new tasks quickly, and reinforcement learning (RL) agents can do the same by transferring knowledge from source policies to a related target task. Transfer RL methods can reshape the policy optimization objective (optimization transfer) or influence the behavior policy (behavior transfer) using source policies. However, selecting the appropriate source policy with limited samples to guide target policy learning has been a challenge. Previous methods introduce additional components, such as hierarchical policies or estimations of source policies' value functions, which can lead to non-stationary policy optimization or heavy sampling costs, diminishing transfer effectiveness. To address this challenge, we propose a novel transfer RL method that selects the source policy without training extra components. Our method utilizes the Q function in the actor-critic framework to guide policy selection, choosing the source policy with the largest one-step improvement over the current target policy. We integrate optimization transfer and behavior transfer (IOB) by regularizing the learned policy to mimic the guidance policy and combining them as the behavior policy. This integration significantly enhances transfer effectiveness, surpasses state-of-the-art transfer RL baselines

in benchmark tasks, and improves final performance and knowledge transferability in continual learning scenarios. Additionally, we show that our optimization transfer technique is guaranteed to improve target policy learning.

Keywords: Optimization transfer, behavior transfer, multi-policy reuse, reinforcement learning

1 Introduction

Trough transferring knowledge from previous policies, humans can learn to solve related new tasks quickly [1]. However, current deep reinforcement learning (RL) agents lack this knowledge transfer ability [2–4], which results in inefficient learning. To address this problem, a large number of research works investigate the multi-policy reuse problem in deep RL: how to efficiently reuse the knowledge from multiple source policies to speed up the learning in a target task [5–9].

To achieve efficient knowledge transfer in RL, the first problem is how to use the source knowledge to influence the learning process in the target task. As there are two major parts in RL: collecting samples and optimizing policies with the collected samples, previous transfer RL works improve the learning efficiency in the target task by either utilizing the source policies to affect the behavior policy of the agent [5, 7, 10], which we name as *behavior transfer*, or reusing the source policies to shape the optimization objective of the target policy [6, 9, 11], which we name as *optimization transfer*. Conducting behavior transfer and optimization transfer is challenging, since there are multiple source policies in the given policy set, and a proper source policy needs to be selected from this set to guide the target policy learning at an early learning stage.

Existing research works learn to select source policies by introducing additional components, such as hierarchical high-level policies over the source policies [7, 8, 10], or estimating the value functions of the source policies on the target task [6, 12, 13]. However, training these additional components significantly harms the transfer effectiveness, as hierarchical policy structures induce a non-stationary issue for policy optimization [14], and estimating the value functions for each source policy is with high sampling cost and computationally expensive. To accomplish efficient transfer without training any additional components, we propose a novel transfer RL method, which employs the value function in the actor-critic framework [15–17] to select the guidance policy from the source policy set, and then uses the selected guidance policy to conduct the transfer. The proposed approach Integrates Optimization transfer and Behavior transfer, which is dubbed as IOB. By inferring the Q function, IOB chooses the source policy that has the largest one-step improvement over the currently learned target policy as the guidance policy. In the policy optimization process, IOB regularizes the target policy to imitate the guidance policy. During the interaction with the environment, the guidance policy and the learned target policy are combined together to form a behavior policy to enable more efficient data collection.

The advantages of the IOB approach are as follows. (i) The one-step improvement can be estimated by querying the Q function and no additional components are needed to be trained. (ii) IOB seamlessly combines optimization transfer and behavior transfer, which accelerates the learning in the target task to a maximum extent. (iii) IOB is conceptually simple and easy to implement, as it introduces very few hyper-parameters to the backbone algorithm. (iv) The optimization transfer in IOB is theoretically guaranteed to improve the target policy learning process. (v) IOB can be naturally integrated with existing continual RL methods to efficiently construct agents with multi-task ability.

To evaluate the proposed method, we compare it with state-of-the-art transfer RL methods on the Meta-World benchmark [18]. Experiment results demonstrate that our method significantly outperforms the baseline methods and achieves the largest *forward transfer*. Then, we visualize the guidance policy selection process to explain the reason why the proposed method works. Next, we perform several ablation studies to analyze the influence of the components of IOB on transfer performance. Finally, we demonstrate that the proposed method could be applied to a continual learning setting, where we combine IOB with a continual learning approach. Experiment results show that IOB boosts the transfer performance of the backbone continual learning approach while maintaining its stability.

We note that a shorter conference version of this paper appeared in [9]. Our initial conference paper has not introduced behavior transfer. This manuscript further promotes the transfer ability of the agent by integrating behavior transfer and optimization transfer. As the guidance policy selection relies on the Q function, this manuscript proposes to increase the accuracy of Q function with ensemble learning. Furthermore, we add a continual learning experiment to better demonstrate the scope of the proposed method.

In the remainder of this paper, we start by reviewing the background knowledge and describing the problem formulation. After that, we review the related work about multi-policy reuse. Next, we present the proposed approach followed by experiment results comparing our approach with the state-of-the-art baselines. Finally, we conclude and outline the directions for future research.

2 Related Works

The learning inefficiency of deep RL approaches restricts their applications to more real-world problems, and transfer learning methods have long been recognized as an effective way to improve the efficiency of the deep RL approaches [19–22]. Here we roughly classify deep transfer RL algorithms into three categories. The algorithms in the first category mainly utilize the source knowledge to reshape the optimization objectives of the benchmark RL methods, which we refer to as the *optimization transfer* methods. The second category focuses on transferring the behavior of the source policies to facilitate the exploration process in the target task, which we refer to as *behavior transfer* methods. The last category is devoted to transferring the parameters of the source policy networks to target policy learning, which we call as *parameter transfer*

methods. These three categories of methods aim to solve the policy optimization, data collection and parameter initialization challenges in deep RL respectively.

Optimization transfer. The optimization transfer methods employ source knowledge to accelerate the optimization process of the target policy. The AC-Teach method [11] uses the value estimation of the source policies to improve the target policy optimization in a Bayesian manner. Similarly, [12] and [6] propose to aggregate the source policies by choosing the policy with the largest Q value at each state. These two methods assume that the source and target tasks share the same dynamics, so that they use the successor features [23] to mitigate the computation cost brought by estimating the Q functions for all the source policies. In contrast, our method only estimates the Q function of the target policy, which is more computationally efficient. Besides, our method works in a more general setting, allowing different dynamic functions for the source and target. The MAMBA method [13] forms a new baseline function by aggregating the value functions of the source policies, and then guides the policy search by improving the policy over the baseline function. The MULTIPOLAR method [24] learns a weighted sum over the actions of the source policies, and learns an auxiliary network to predict the residuals around the aggregated actions. Compared with these previous methods, our method does not require training any additional components, which is both computationally efficient and sampling efficient.

Behavior transfer. The behavior transfer methods aim at improving the exploration efficiency in the target task with the given source policies. A series of works propose to combine the source policies with random policies probabilistically to achieve more efficient exploration [5, 25, 26]. As the combination manner is not learned, these methods cannot guarantee to perform better than learning without the source knowledge. To accomplish a more effective policy combination, the following works [7, 8, 10, 27] propose a hierarchical policy structure to reuse the given policies, where a high-level policy is learned to select which source policy should be executed at the current state. Although the policy combination is more accurate under the hierarchical structure, the simultaneous learning of multi-level policies suffer from the non-stationary issue [14]. In this work, we conduct behavior transfer without the hierarchical policy structure. Instead, the source policies and the target policy are combined under the guidance of the learned critic.

Parameter transfer. The parameter transfer methods initialize the neural networks for the target task with the parameters learned in the source tasks. When there are multiple source tasks, it is a challenging problem which source policy to transfer from, since the parameter initialization needs to be conducted before the learning in the new target task starts. Some works propose a progressive network structure which connects these source neural networks with lateral connections [28, 29], and then use the progressive network as the initialization in the target task. These methods may not be scalable when the number of the source tasks grow too large. For better scalability, some following works propose to prune the source networks or distill the source networks to a smaller network and then reuse the parameters [30–32]. Note that our method and the parameter transfer methods are orthogonal, and could be easily combined together. In section 5.5, we combine the proposed method with a parameter

transfer method based on pruning called PackNet [31], and evaluate it in a continual learning setting.

3 Preliminaries and Problem Formulation

The environment in RL is formulated as a Markov Decision Process (MDP), and the MDP is defined by a tuple $(\mathcal{S}, \mathcal{A}, p, r, \gamma)$, where \mathcal{S} is a state space, \mathcal{A} is an action space, $p(s'|s, a)$ is an unknown transition function, $r(s, a, s')$ is a reward function, and $\gamma \in [0, 1)$ is a discount factor. The objective of RL is to learn a policy $\pi(a|s)$ that could maximize the expected discounted return: $R(\pi) = \mathbb{E}_\pi [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, s_{t+1}) | a_t \sim \pi(\cdot|s), s_0]$.

While the proposed approach could be easily integrated with off-policy actor-critic methods, in the following sections, we mainly present how it could be combined with the Soft Actor-Critic (SAC) algorithm [17]. By removing the entropy term in SAC, the proposed approach can be applied to other off-policy actor-critic RL methods as well, such as Deep Deterministic Policy Gradient (DDPG) [15] and Twin-Delayed DDPG (TD3) [16]. Next, we introduce some preliminary knowledge about the SAC method.

SAC: The SAC method [33] introduced an additional function approximator for the value function, but later found it to be unnecessary [17]. In this paper, the soft Q function and soft V function of a policy π in SAC are defined as:

$$\begin{aligned} Q_\pi(s, a) &= r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot|s, a)} [V_\pi(s)], \\ V_\pi(s) &= \mathbb{E}_{a \sim \pi(\cdot|s)} [Q_\pi(s, a) - \alpha \log \pi(a|s)], \end{aligned} \quad (1)$$

where $\alpha > 0$ is the entropy weight, and the loss functions of SAC are defined as:

$$\begin{aligned} L_{critic}(Q_\theta) &= \mathbb{E}_{(s, a, r, s') \sim \mathcal{D}} [(Q_\theta(s, a) - (r + \gamma V_{\bar{\theta}}(s')))]^2, \\ L_{actor}(\pi_\phi) &= \mathbb{E}_{s \sim \mathcal{D}} [\mathbb{E}_{a \sim \pi_\phi(\cdot|s)} [\alpha \log \pi_\phi(a|s) - Q_\theta(s, a)]], \\ L_{entropy}(\alpha) &= \mathbb{E}_{s \sim \mathcal{D}} [\mathbb{E}_{a \sim \pi_\phi(\cdot|s)} [-\alpha \log \pi_\phi(a|s) - \alpha \bar{\mathcal{H}}]], \end{aligned} \quad (2)$$

where \mathcal{D} is the replay buffer, $\bar{\mathcal{H}}$ is a hyper-parameter representing the target entropy, θ and ϕ are network parameters, $\bar{\theta}$ denote the parameters of the target network, and $V_{\bar{\theta}}(s) = \mathbb{E}_{a \sim \pi(a|s)} [Q_{\bar{\theta}}(s, a) - \alpha \log \pi(a|s)]$ is the target soft value function.

Based on the SAC method, we define the *soft expected advantage* of action probability distribution $\pi_i(\cdot|s)$ over policy π_j at state s as:

$$Adv_{\pi_j}(s, \pi_i) = \mathbb{E}_{a \sim \pi_i(\cdot|s)} [Q_{\pi_j}(s, a) - \alpha \log \pi_i(a|s) - V_{\pi_j}(s)]. \quad (3)$$

$Adv_{\pi_j}(s, \pi_i)$ measures the one-step performance improvement brought by following π_i instead of π_j at state s , and following π_j afterwards.

Problem Formulation: Multi-policy reuse focuses on learning the policy for a target MDP M with fewer samples through transferring knowledge from a set of source policies $\{\pi_1, \pi_2, \dots, \pi_n\}$. We denote the target policy learned on M as π_{tar} , and its corresponding soft Q function as $Q_{\pi_{tar}}$. In this paper, we assume that the source policies and the target policy share the same state and action spaces. This assumption

generally holds for the tasks with homogeneous agents, e.g., one robot manipulating different objects, or one robot with the same reception field navigating in different environments.

4 Method

There are two prominent components that significantly affect the efficiency of RL: the policy optimization objective and the way of collecting samples. Although the goal of most RL methods is to maximize the expected return, it remains an unsolved problem how to use the source knowledge to shape the policy optimization objective, so that achieving this goal costs fewer stochastic gradient descent iterations. Beyond that, the way of collecting samples (behavior policy) also plays a crucial role to improve learning efficiency, as the behavior policy determines the quality of the training data of the neural networks. In this paper, we propose a novel transfer RL approach, which aims at improving the learning efficiency in the target task from both the optimization and behavior perspectives using a source policy set.

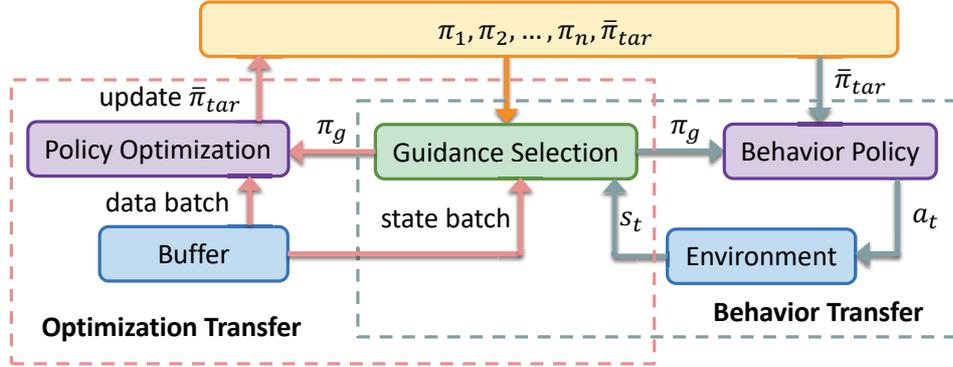


Fig. 1 The overall framework of the proposed approach.

Figure 1 visualizes the overall framework of the proposed approach, Integrating Optimization transfer and Behavior transfer for multi-policy reuse (IOB). The left dash box contains the flowchart of *optimization transfer*, which is described in Section 4.1, and the right dash box illustrates the *behavior transfer* process, which is presented in Section 4.2. Then, in Section 4.3, we elaborate on how those two types of transfer are combined together to thoroughly boost the target task learning efficiency. Furthermore, in Section 4.4, we analyze the proposed method from a theoretical perspective and prove that even with an approximated Q function, the target policy is guaranteed to be improved monotonically with the proposed optimization transfer technique. Finally, in Section 4.5, We have integrated IOB with the continual learning method to propose a novel continual RL approach highlighting transfer capabilities.

4.1 Optimization Transfer

To achieve positive transfer from the optimization perspective, IOB utilizes action distributions output by the source policies to guide the learning of the target policy. Specifically, at state s , the agent has access to a set of candidate action distributions output by $n + 1$ policies, including n source policies and a hard copy of the target policy, $\bar{\pi}_{tar}$:

$$\Pi^s = \{\pi_1(\cdot|s), \pi_2(\cdot|s), \dots, \pi_n(\cdot|s), \bar{\pi}_{tar}(\cdot|s)\}. \quad (4)$$

From this candidate set, IOB selects a guidance policy π_g which has the largest one-step improvement over the current target policy:

$$\begin{aligned} \pi_g(\cdot|s) &= \arg \max_{\pi(\cdot|s) \in \Pi^s} Adv_{\pi_{tar}}(s, \pi) \\ &= \arg \max_{\pi(\cdot|s) \in \Pi^s} \mathbb{E}_{a \sim \pi(\cdot|s)} [Q_{\pi_{tar}}(s, a) - \alpha \log \pi(a|s)]. \end{aligned} \quad (5)$$

The second equation holds as adding $V_{\pi_{tar}}(s)$, to all soft expected advantages does not affect the result of the $\arg \max$ operator. Note that the guidance policy selection is conditioned on states, and Equation (5) indicates that the action output by the guidance policy π_g at state s is at least no worse than the current target policy π_{tar} in terms of the expected return estimated by the Q value. Possibly the guidance policy is better than the current target policy, if the source policy set contains one or more policies similar to the optimal target policy. As we have obtained a guidance policy, the next question is how we could use it to guide the target policy update. Since the guidance policy π_g may accomplish larger returns than the current target policy, we propose to regularize the target policy π_{tar} to imitate the guidance policy π_g selected from the candidate set Π^s before this update, and minimize the following loss function:

$$\begin{aligned} L_{\pi}(\pi_{tar}) &= L_{actor}(\pi_{tar}) \\ &\quad + \mathbb{E}_{s \sim \mathcal{D}} [\beta_s D_{KL}(\pi_{tar}(\cdot|s) || \pi_g(\cdot|s))], \end{aligned} \quad (6)$$

where L_{actor} is the original actor loss defined in Equation (2), and $\beta_s > 0$ is a hyper-parameter controlling the weight of the regularization. The training data for the regularization is also sampled from the replay buffer D , the same as that of L_{actor} . After each update of the target policy π_{tar} , the corresponding policy in the candidate policy set is synchronized immediately.

Since the Q value is critical to the guidance policy selection in Equation (5), we need a relatively accurate estimation of the Q value, so that the guidance policy could be beneficial to the target policy optimization. As the value function learning in RL often suffers from the over-estimation issue [34, 35], we propose to apply the following critic ensemble technique to attain a more accurate Q function:

$$Q_{\pi_{tar}}(s, a) = \min_{k \in [1..K]} Q_{\theta_k}(s, a), \quad (7)$$

where θ_k denotes the parameters of the k -th Q network. All the Q networks are independently initialized and trained. Limited to the computation resource, also balancing

overestimation and underestimation of those Q functions, the total number for the Q-networks, K , is set as 4 in the experiment section.

4.2 Behavior Transfer

The regularization in Section 4.1 enables faster policy learning with the training data sampled from the replay buffer, and another problem is how to fill the replay buffer with high-quality data with large returns. To solve this challenge, we propose to further employ the guidance policy to improve the behavior policy. As most off-policy algorithms could only afford a slight degree of *off-policyness* [36, 37], the proposed approach probabilistically combines the guidance policy and the learned target policy as the behavior policy in an ϵ -greedy manner.

Algorithm 1 Behavior- $\pi(\epsilon, s_t, \Pi_s, \bar{\pi}_{tar}, \{Q_{\theta_k} | k \in 1..K\})$

```

1: if  $random() < \epsilon$  then
2:    $\pi_b \leftarrow \arg \max_{\pi(\cdot|s_t) \in \Pi^s} \min_{k \in 1..K} Q_{\theta_k}(s_t, a) - \alpha \log \pi(a|s_t),$ 
3:   where  $a \sim \pi(\cdot|s_t)$ 
4: else
5:    $\pi_b \leftarrow \bar{\pi}_{tar}$ 
6: end if
7:  $a_t \sim \pi_b(\cdot|s_t)$ 
8: return  $a_t$ 

```

As $\bar{\pi}_{tar}$ is synchronized with π_{tar} immediately after each policy update, the output of $\bar{\pi}_{tar}$ is the same as π_{tar} , so we use $\bar{\pi}_{tar}$ as the current target policy in Algorithm 1. At the beginning of the learning process, the Q-value estimation may be inaccurate, and this behavior policy could be regarded as the optimistic exploration towards the actions with overestimated Q values. To limit the off-policyness of the behavior policy, ϵ in Algorithm 1 needs to be a small value. However, when ϵ approaches 0 too much, the behavior policy cannot take advantage of the guidance policy. To balance exploration and exploitation, we set $\epsilon = 0.2$ in the experiments.

4.3 Integrating Optimization Transfer and Behavior Transfer

The pseudo-code of IOB is presented in Algorithm 2. When interacting with the environment, the agent probabilistically utilizes the guidance policy to collect samples (Line:6-11). During the policy updates, the guidance policy regularizes the direction of target policy optimization to achieve more efficient learning (Line:15-18). After each policy update, $\bar{\pi}_{tar}$ in the policy set Π_s is synchronized from the target policy π_{tar} (Line 20). Furthermore, as the guidance selection heavily depends on the learned critic, to select an effective guidance policy, we need a well-estimated Q function. Therefore, we employ the critic ensemble technique, and use different data to update multiple Q networks (Line:13-14).

Algorithm 2 IOB

```
1: Require: Source policies  $\{\pi_1, \pi_2, \dots, \pi_n\}$ , hyper-parameters  $\lambda_\pi, \lambda_\alpha, \tau, \overline{\mathcal{H}}, \beta_s, \epsilon$ 
2: Initialize replay buffer  $\mathcal{D}$ 
3: Initialize  $\pi_{tar}$  with parameter  $\phi$ , entropy weight  $\alpha$ , critic  $Q_{\theta_k}$ , target critic  $Q_{\overline{\theta}_k} \leftarrow Q_{\theta_k}$ , for  $k \in \{1..K\}$ 
4:  $\overline{\pi}_{tar} \leftarrow \pi_{tar}$ ,  $\Pi^s \leftarrow \{\pi_1(\cdot|s), \pi_2(\cdot|s), \dots, \pi_n(\cdot|s), \overline{\pi}_{tar}(\cdot|s)\}$ 
5: while not done do
6:   for each environment step do
7:      $a_t \leftarrow \text{Behavior-}\pi(\epsilon, s_t, \Pi_s, \overline{\pi}_{tar}, \{Q_{\theta_k} | k \in 1..K\})$ 
8:      $s_{t+1} \sim p(s_{t+1}|s_t, a_t)$ 
9:      $\mathcal{D} \leftarrow \mathcal{D} \cup \{s_t, a_t, r(s_t, a_t), s_{t+1}\}$ 
10:     $s_t \leftarrow s_{t+1}$ 
11:   end for
12:   for each gradient step do
13:     Sample  $K$  minibatches from  $\mathcal{D}$ , and update the critic networks independently
14:      $\overline{\theta}_k \leftarrow \tau\theta_k + (1 - \tau)\theta_k$  for  $i \in \{1..K\}$ 
15:     Sample minibatch  $b$  from  $\mathcal{D}$  to update  $\pi_{tar}$  and  $\alpha$ 
16:     Query the action probabilities  $\{\pi_1(\cdot|s), \pi_2(\cdot|s), \dots, \pi_n(\cdot|s), \overline{\pi}_{tar}(\cdot|s)\}$  for state  $s$  in  $b$ 
17:     Compute expected advantages according to Eq. (3), form  $\pi_g$  according to Eq. (5)
18:      $\phi \leftarrow \phi - \lambda_\pi \hat{\nabla}_\phi L_\pi(\pi_{tar})$ 
19:      $\alpha \leftarrow \alpha - \lambda_\alpha \hat{\nabla}_\alpha L_{entropy}(\alpha)$ 
20:     Synchronize  $\Pi_s$  with the updated  $\pi_{tar}$ 
21:   end for
22: end while
23: return  $\pi_{tar}$ 
```

4.4 Theoretical Analysis

Note that we can hardly acquire the exact Q values to select the guidance policy during learning, and the Q values are estimated with function approximation in deep RL. In this subsection, we provide a theoretical analysis that even with an approximated Q function, we could form the guidance policy, and by regularizing the target policy to mimic the guidance policy, the target policy learning is guaranteed to achieve a monotonic improvement.

Theorem 1. *Let $\tilde{Q}_{\pi_{tar}}$ be an approximation of $Q_{\pi_{tar}}$, s.t.,*

$$|\tilde{Q}_{\pi_{tar}}(s, a) - Q_{\pi_{tar}}(s, a)| \leq \mu \text{ for all } s \in \mathcal{S}, a \in A. \quad (8)$$

Define

$$\tilde{\pi}_g(\cdot|s) = \arg \max_{\pi(\cdot|s) \in \Pi^s} \mathbb{E}_{a \sim \pi(\cdot|s)} \left[\tilde{Q}_{\pi_{tar}}(s, a) - \alpha \log \pi(a|s) \right], \quad (9)$$

for all $s \in \mathcal{S}$.

Then,

$$V_{\tilde{\pi}_g}(s) \geq V_{\pi_{tar}}(s) - \frac{2\mu}{1-\gamma} \text{ for all } s \in \mathcal{S}. \quad (10)$$

Theorem 1 provides a way to obtain the guidance policy with approximated Q values, and the SAC method naturally learns that approximation, so that the guidance policy could be formed without training any additional components. In the following, we provide another theorem, which proves that policy improvement can be guaranteed if the target policy is optimized to stay close to the guidance policy.

Theorem 2. *If*

$$D_{KL}(\pi_{tar}^{l+1}(\cdot|s) || \tilde{\pi}_g^l(\cdot|s)) \leq \delta \text{ for all } s \in \mathcal{S}, \quad (11)$$

then

$$V_{\pi_{tar}^{l+1}}(s) \geq V_{\pi_{tar}^l}(s) - \frac{\sqrt{2 \ln 2\delta}(\tilde{R}_{max} + \alpha \mathcal{H}_{max}^{l+1})}{(1-\gamma)^2} - \frac{2\mu + \alpha \tilde{\mathcal{H}}_{max}}{1-\gamma} \text{ for all } s \in \mathcal{S}, \quad (12)$$

where $\tilde{\pi}_g^l$ is the guidance policy selected after the l -th policy update, π_{tar}^l and π_{tar}^{l+1} is the learned target policy after the l -th and $(l+1)$ -th policy update. $\tilde{R}_{max} = \max_{s,a} |r(s,a)|$ is the largest possible absolute value of the reward, $\mathcal{H}_{max}^{l+1} = \max_s \mathcal{H}(\pi_{tar}^{l+1}(\cdot|s))$ is the largest entropy of π_{tar}^{l+1} , and $\tilde{\mathcal{H}}_{max} = \max_s |\mathcal{H}(\pi_{tar}^l(\cdot|s)) - \mathcal{H}(\pi_{tar}^{l+1}(\cdot|s))|$ is the largest possible absolute difference of the policy entropy.

4.5 Transfer in a Continual RL Setting

In Continual RL [38], an agent will sequentially learn a series of tasks $\mathcal{Z}^{(1)}, \dots, \mathcal{Z}^{(T_{max})}$, each corresponding to an individual MDP $\mathcal{Z}^{(t)} = (\mathcal{S}^{(t)}, \mathcal{A}^{(t)}, p^{(t)}, r^{(t)}, \gamma)$, while maintaining fixed constraints on computation and memory. The agent seeks out an optimal set of policy parameters $\{\phi^{(1)}, \dots, \phi^{(T_{max})}\}$ to maximize the average rewards across all tasks. The Continual RL algorithm requires stability, i.e., the ability to prevent forgetting acquired skills, and plasticity, i.e., the ability to learn new skills quickly. Measures that enhance only one of these abilities often limit the other, resulting in a stability-plasticity dilemma.

To balance stability and plasticity as well as boost transfer, we combine the proposed method IOB with an advanced continual learning method, PackNet [31]. [39] compared seven representative continual RL methods under the sequence of robotic arm tasks and showed that PackNet outperformed all other methods. PackNet develops a training-pruning-retraining procedure. After pruning, the parameters belonging to the previously learned policy are frozen, and only the pruned parameters could be updated in the following tasks, so that the PackNet approach hardly forgets any policy.

However, the transfer ability of PackNet is limited, as it only considers parameter transfer by representing all the policies with one neural network. Integrating IOB with PackNet can potentially enhance the agent's ability to build upon prior knowledge. We wonder if IOB could improve the transfer performance of PackNet with optimization

transfer and behavior transfer while maintaining the property of no forgetting. When applying IOB in the continual learning setting, we treat all the previously learned policies as source policies, i.e., when learning in the i -th task, there are $i - 1$ source policies. The pseudocode of integrating IOB with PackNet is presented in Algorithm 3.

Algorithm 3 Continual RL with IOB

- 1: **Require:** Continuous reinforcement learning task sequences $\mathcal{Z}^{(1)}, \dots, \mathcal{Z}^{(T_{max})}$, hyper-parameters $\lambda_\pi, \lambda_\alpha, \tau, \overline{\mathcal{H}}, \beta_s, \epsilon$
 - 2: Initialize source policies $\Pi^s = \{\}$
 - 3: **for** each task $\mathcal{Z}^{(t)}$ **do**
 - 4: $\pi_t = \text{IOB}(\Pi^s, \text{hyper-parameters } \lambda_\pi, \lambda_\alpha, \tau, \overline{\mathcal{H}}, \beta_s, \epsilon)$
 - 5: Sort parameters in each layer
 - 6: Prune the smallest $\frac{n-t}{n-t+1}\%$ parameters per layer
 - 7: Retrain network to maintain performance on $\mathcal{Z}^{(t)}$
 - 8: Freeze parameters for $\mathcal{Z}^{(t)}$ and never update them
 - 9: Append π_t to Π^s
 - 10: **end for**
 - 11: **return** Π^s
-

5 Experiments

We evaluate the proposed method IOB on Meta-World [18], a popular RL benchmark composed of multiple robotic manipulation tasks. These tasks are both correlated (performed by the same Sawyer robot arm) and distinct (interacting with different objects and having different reward functions), and therefore serve as a proper evaluation benchmark for policy reuse. The source policies are obtained by training on three representative tasks: Reach, Push, and Pick-Place. We choose several complex tasks as target tasks, including Hammer, Peg-Insert-Side, Push-Wall, Pick-Place-Wall, Push-Back, and Shelf-Place. Among these target tasks, Hammer and Peg-Insert-Side require interacting with objects unseen in the source tasks. In Push-Wall and Pick-Place-Wall, there is a wall between the object and the goal. In Push-Back, the goal distribution is different from Push. In Shelf-Place, the robot is required to put a block on a shelf, and the shelf is unseen in the source tasks. Figure 2 visualizes these tasks and the video demonstrations are available at <https://meta-world.github.io/>. Similar to the settings in [40], in our experiments the goal position is randomly reset at the start of every episode, which increases the stochasticity in the environment and is thus more challenging.

In this section, we first briefly introduce the baseline methods and the implementation details. Next, we show the comparison results in the Meta-World benchmark. Then, we analyze the guidance policy selection to dive into the reason why the proposed method could achieve positive transfer. After that, we perform several ablation studies to analyze the influence of the components of IOB on transfer performance.

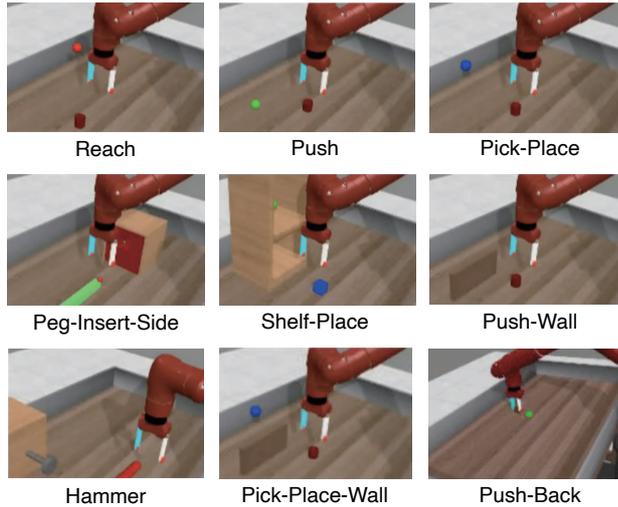


Fig. 2 Visualization of the tasks in the experiment section.

Finally, we conduct a continual learning experiment, which demonstrates the effectiveness of the proposed method when the agent continuously learns in a sequence of tasks.

5.1 Baselines and Implementation Details

We compare the proposed method with several representative transfer RL baseline algorithms and the backbone RL method in this work.

- CUP: A critic-guided policy reuse method without behavior transfer [9].
- HAAR: A hierarchical policy reuse method that simultaneously learns two-level policies [7].
- MAMBA: An optimization transfer method based on value function aggregation [13].
- MULTIPOLAR; A transfer RL method learning a weighted sum of the source policies’ action probabilities and an additional network to predict residuals [24].
- SAC: An off-policy actor-critic method with entropy maximization [17].

Implementation details: To improve the computation efficiency of IOB, we store the outputs of the source policies in the replay buffer. When saving the buffer, we can query those source policies with batches of states. As those outputs are stored in the replay buffer, each state only needs to be queried one time, which leads to better computation efficiency.

Equation (5) requires estimating the expectation over Q values. In practice, we obtain the estimation by sampling a few actions (5 actions) from the action distributions output by the source policies, and find that it is sufficient to accomplish a stable performance. The hyper-parameters used in the experiments are listed in Table 1. The

same set of hyper-parameters is used for all the tasks, and most hyper-parameters are adopted from [41].

Table 1 Hyper-parameter settings.

Hyper-Parameter	Value
batch size	1280
non-linearity	ReLU
actor network structure	3 fully-connected layers with 400 units
critic network structure	3 fully-connected layers with 400 units
policy initialization	standard Gaussian
learning rates for all networks	3e-4
optimizer	Adam
episode length (horizon)	500
discount	0.99
regularization rate β_s	30
ensemble number K	4
guidance policy prob ϵ	0.2

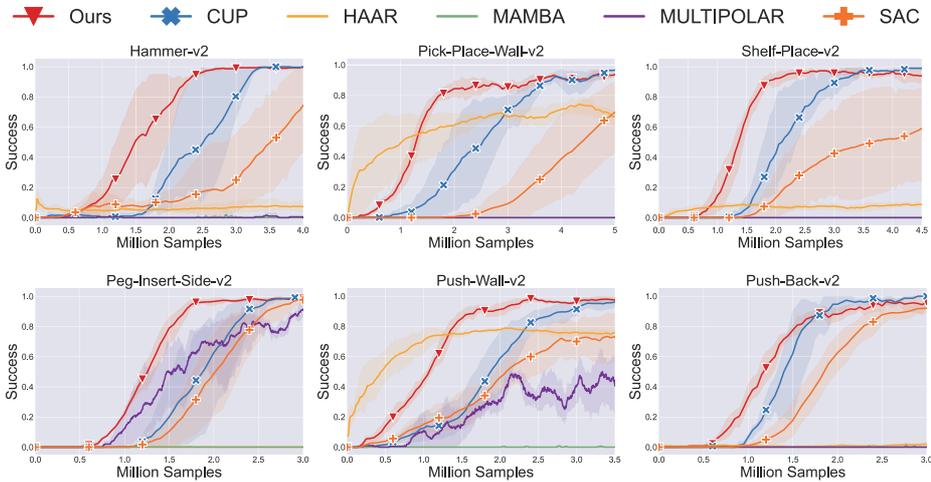


Fig. 3 Learning curves of the proposed method and the baselines on the tasks of the Meta-World benchmark. Videos comparing the policy learning processes of our method and SAC are at <https://sites.google.com/view/iob-aamas>.

5.2 Comparison Results on Meta-World

We present the learning curves of all the methods in Figure 3. Those learning curves are averaged over 5 runs, and the success rates are averaged by 10 evaluations. We also compare the *forward transfer* of these methods in Table 2. Forward transfer, FT , is quantitative measure in transfer RL [39], which is defined as the normalized area

between the learning curve of the transfer RL method and the learning curve of the reference method without transfer. Specifically, we use SAC as the reference method. Figure 4 is an example of forward transfer.

$$FT = \frac{AUC_{trans} - AUC_{SAC}}{1 - AUC_{SAC}}. \quad (13)$$

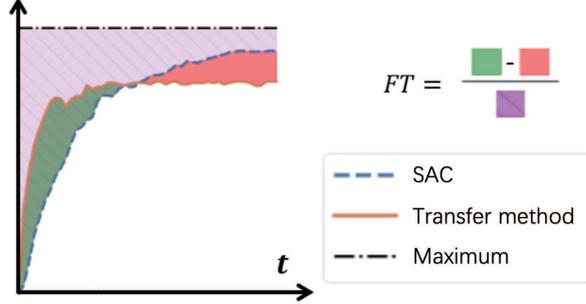


Fig. 4 An example of forward transfer FT [39].

Table 2 Forward Transfer Compared to SAC

	Ours	CUP	HAAR	MAMBA	MULTIPOLAR
Hammer	0.51 ± 0.09	0.25 ± 0.10	-0.15 ± 0.02	-0.23 ± 0.0	-0.23 ± 0.0
Peg-Insert-Side	0.37 ± 0.02	0.08 ± 0.13	-0.45 ± 0.0	-0.45 ± 0.0	0.15 ± 0.14
Pick-Place-Wall	0.61 ± 0.01	0.38 ± 0.12	0.59 ± 0.02	-0.18 ± 0.0	-0.18 ± 0.0
Push-Wall	0.49 ± 0.04	0.13 ± 0.12	0.58 ± 0.02	-0.58 ± 0.0	-0.22 ± 0.08
Shelf-Place	0.56 ± 0.03	0.35 ± 0.07	-0.21 ± 0.01	-0.31 ± 0.0	-0.31 ± 0.0
Push-Back	0.34 ± 0.04	0.26 ± 0.07	-0.49 ± 0.0	-0.51 ± 0.0	-0.51 ± 0.0
Avg	0.48	0.24	-0.02	-0.38	-0.22

As shown in Figure 3, our approach has achieved significantly better performance than the baseline methods. As lacking behavior transfer, the CUP method is less efficient than the proposed approach. Note that the guidance policy selection for IOB and CUP is the same, comparing with CUP could be regarded as an ablation study for behavior transfer. HAAR has a jump-start performance on Push-Wall and Pick-Place-Wall. However, due to the non-stationary issue induced by jointly training the high-level and low-level policies in the hierarchical structure, HAAR can hardly converge to a success rate close to 1. MULTIPOLAR achieves better performance on Push-Wall and Peg-Insert-Side than on the other tasks, because the Push source policy is useful on Push-Wall (implied by HAAR’s good jump-start performance), and learning residuals on Peg-Insert-Side is easier (implied by SAC’s fast learning). In Pick-Place-Wall, the Pick-Place source policy is useful, but the residual is challenging to learn, so MULTIPOLAR does not work. For the remaining three tasks, the source policies are less useful, and therefore MULTIPOLAR fails on these tasks. The MAMBA method

could hardly learn any successful policy, since accurately estimating the value functions for all source policies is not sample efficient.

Table 2 demonstrates the forward transfer values for all the methods. In almost all the tasks, the proposed approach has achieved the largest forward transfer, which indicates that our approach generally works. In the Push-Wall task, the forward transfer of HAAR is slightly larger than ours, but this slight outperformance is at the expense of a smaller convergent success rate, which is unfavorable.

5.3 Analysis of Guidance Policy Selection

This subsection visualizes the guidance policy selection during the learning process of the proposed approach. Figure 5 shows the percentages of the given source policies and the learned target policy being selected as the guidance policy throughout the training on the Push-Back task. Note that the critic network is randomly initialized, and after several optimization iterations, the critic could identify which action is more beneficial. Therefore, after 0.25M steps, we utilize the guidance policy selected with the critic to conduct transfer.

At the early stages of training, the source policies are selected more frequently as they have positive expected advantages. This indicates that they can be used to improve the current target policy. As the training proceeds and the target policy becomes better, the source policies are selected less frequently. Among these three source policies, Reach is chosen more frequently than the other two source policies, as in most manipulation tasks, the robot needs to first reach the target object. Figure 5 validates that the proposed method could transfer knowledge from multiple source policies to facilitate the target task learning as well.

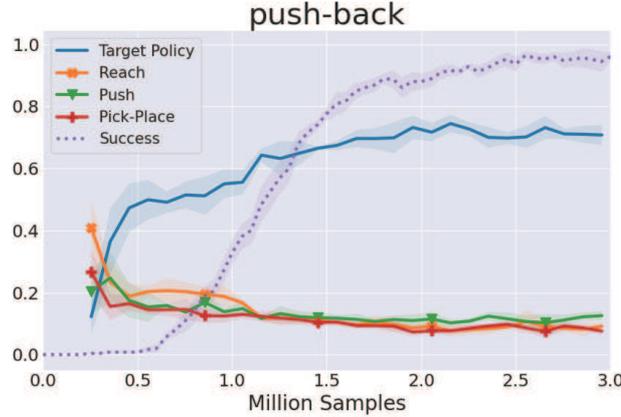


Fig. 5 Percentages of the source policies and the current target policy selected as the guidance policy during training in the Push-Back task. The dashed line is the success rate of the target policy in this task.

5.4 Ablation Study

In this subsection, we conduct three ablation studies to answer the following questions: (1) Could the proposed method benefit from a larger source policy set, which includes more related source policies? (2) What is the influence of random source policies on the transfer performance of the proposed method? (3) Could behavioral transfer enhance transfer performance, and what is the proper value of ϵ in the evaluated tasks?

5.4.1 Larger Source Policy Sets

We evaluate the proposed method with a larger source policy set on the Pick-Place-Wall task. The original source policy set is expanded with three additional policies, which solve the Drawer-Close, Push-Wall and Coffee-Button tasks, i.e., the new source policy set is composed of 6 policies. Figure 6 provides the comparison results of our method with 3 source policies and with 6 source policies. The results show that our method is able to utilize the additional source policies to further improve the transfer performance. As the original policy set with three policies has already contained one useful source, Pick-Place, the improvement caused by the additional source policies is relatively slight.

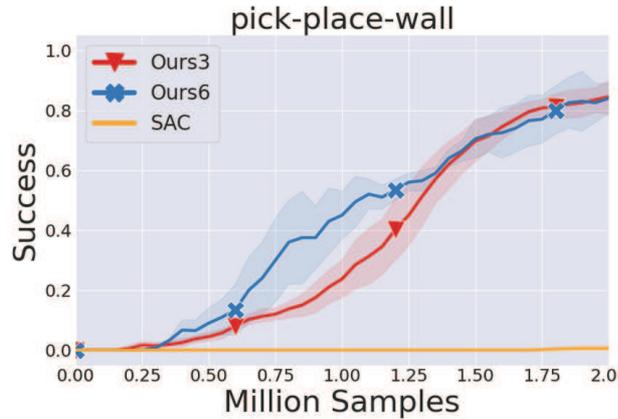


Fig. 6 Ablation study with different numbers of source policies on the Pick-Place-Wall task.

5.4.2 Random Source Policies

To investigate the robustness of the proposed method to the useless random source policies, we design a transfer setting with two source policy sets. One source policy contains three random policies, and the other contains the Reach policy and three random policies. As shown in Figure 8, when no source policy is useful, the proposed method performs similarly to the SAC method, and its sample efficiency is almost unaffected by those random source policies. When there is one useful source policy, the proposed method could efficiently utilize it to improve the learning performance, even if a lot of useless random source policies exist.

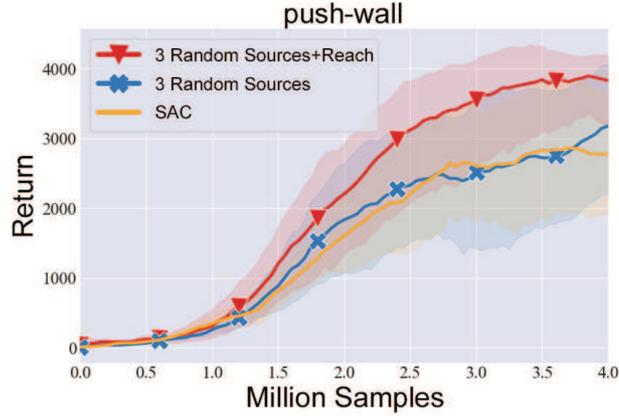


Fig. 7 Ablation studies on IOB’s sensitivity to useless source policies.

5.4.3 Selection of Hyperparameter ϵ

In order to demonstrate the effectiveness of the behavior transfer, we evaluate the transfer performance of IOB under different ϵ in the pick-place-wall-v2 task. We have visualized the early training process of this task in Figure 8. Evidently, without behavior transfer, the agent cannot obtain high-value samples solely through self-exploration. Conversely, appropriate behavior transfer can effectively enhance transfer performance; however, excessive behavioral transfer can result in performance decline. This can be attributed to the model learned almost exclusively under offline settings, making it more challenging to evaluate the Q function accurately. Notably, accurate evaluation of the Q-function is crucial for the IOB method.

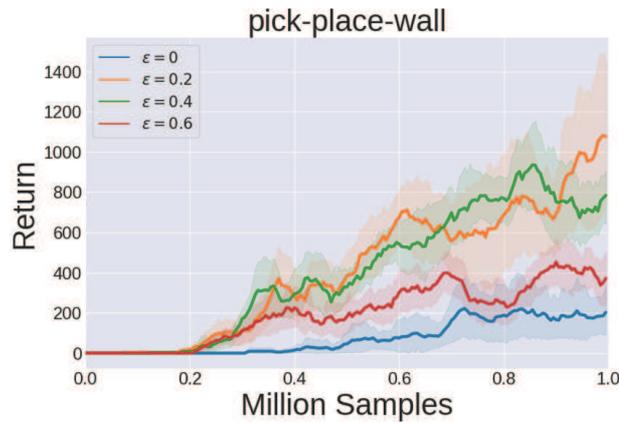


Fig. 8 Ablation study on behavioral transfer effects over a range of ϵ magnitudes.

5.5 Transfer in a Continual Learning Setting

To validate whether the proposed method could enhance transfer in a continual learning setting, we conduct a continual RL experiment. In this experiment, an agent sequentially learns policies in eight tasks: Reach-v2, Push-v2, Pick-Place-v2, Hammer-v2, Shelf-Place-v2, Peg-Insert-Side-v2, Push-Wall-v2, Pick-Place-Wall-v2. In each task, the sample budget is 3 million samples. We compare our method with PackNet and a naive continual learning baseline, i.e., fine-tuning the previously learned policy. As continual learning requires the agent to have the ability of both learning new policies and not forgetting previously learned policies, we evaluate these policies on all the eight tasks, and report the average success rates over tasks in Figure 9.

The results in Figure 9 demonstrate that the PackNet method remember all the previously learned policy, as the average success rate curve of PackNet has not dropped in the learning process. Compared with PackNet, the forgetting phenomenon of fine-tuning is very obvious. At the beginning of learning in a new task (the dashed lines), the average success rate of fine-tuning drops immediately, which indicates that the fine-tuning method forgets the previous policy. Taking advantage of the parameter isolation technique of PackNet, our method could remember all the previously learned policies as well. Beyond that, our method achieves faster learning and larger convergence success rates via behavior transfer and optimization transfer. After training in the 8 tasks with 24 (3×8) million steps, our method has achieved an average success rate of more than 0.8.

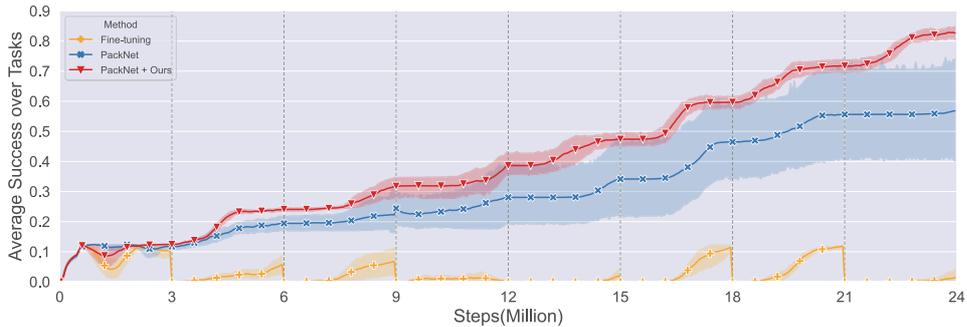


Fig. 9 Average success rates over the 8 tasks of Fine-tuning, PackNet, and applying the proposed method IOB to PackNet.

In Table 3, we provide the average convergent success rate at T steps ($T = 24$ million), average forward transfer, and average forgetting over the 8 tasks. The forgetting of the i -th policy is measured by the drop of the success rate from the end of training in the i -th task to the end of the total learning process,

$$Forgetting_i = SR(i \cdot \Delta) - SR(T), \quad (14)$$

where SR denotes the evaluation success rate, Δ is the sample budget for each task (3 million). From Table 3, we can see that our method has achieved the largest forward

transfer, the least forgetting and the best convergent success rate among the three methods. However, compared to the average forward transfer over the tasks in Table 2, we find that the forward transfer of combining our method with PackNet is smaller, which indicates that the effect of reusing parameters may be negative when conducting optimization transfer and behavior transfer. The PackNet method demonstrates no forgetting, but its forward transfer is smaller than its variant combined with our method.

Table 3 Continual Learning Results

	Success	Forgetting	Transfer
Fine-tuning	0.02 ± 0.02	0.55 ± 0.12	-0.06 ± 0.14
PackNet	0.57 ± 0.20	0.00 ± 0.0	0.07 ± 0.14
PackNet+Ours	0.83 ± 0.02	0.01 ± 0.01	0.24 ± 0.13

6 Conclusion

To transfer knowledge efficiently from multiple source policies to a related target task, we propose a novel transfer RL method that conducts guidance policy selection without training any extra component. By utilizing the Q function as a natural evaluation of the source policies, the proposed method selects the policy with the largest one-step improvement over the current target policy as the guidance policy. The selected guidance policy is used to regularize the policy optimization process to accomplish optimization transfer. Meanwhile, to boost transfer performance further with behavior transfer, we combine the guidance policy with the learned policy as the behavior policy. Benefiting from the effective guidance policy selection and the integration of optimization transfer and behavior transfer, the proposed method significantly outperforms the state-of-the-art transfer RL baselines on the benchmark tasks. Beyond that, we provide a theoretical analysis that with the proposed policy optimization technique, the target policy is guaranteed to be improved monotonically.

As for future work, we consider relaxing the assumption in this paper that all the source policies and the target policy share the same state and action spaces. This assumption somehow limits the applications of the proposed method to more general circumstances. Aligning different state and action spaces is a challenging problem. Previous research works investigate this problem with a high-level structure [42–46], and some inspiration could be taken from those previous works.

Acknowledgments. This work is supported by the Key Program of the National Natural Science Foundation of China (Grant No.51935005), Basic Research Project (Grant No.JCKY20200603C010), Natural Science Foundation of Heilongjiang Province of China (Grant No.LH2021F023), as well as Science and Technology Planning Project of Heilongjiang Province of China (Grant No.GA21C031), and China Academy of Launch Vehicle Technology (CALT2022-18).

Proof of Theorem 1

Proof. As $|\tilde{Q}_{\pi_{tar}}(s, a) - Q_{\pi_{tar}}(s, a)| \leq \mu$ for all $s \in \mathcal{S}, a \in A$, we have that for all $s \in \mathcal{S}$, the difference between the true value function $V_{\pi_{tar}}$ and the approximated value function $\tilde{V}_{\pi_{tar}}$ is bounded:

$$\begin{aligned} & V_{\pi_{tar}}(s) \\ &= \mathbb{E}_{a \sim \pi_{tar}(\cdot|s)} [Q_{\pi_{tar}}(s, a) - \alpha \log \pi_{tar}(a|s)] \\ &\leq \mathbb{E}_{a \sim \pi_{tar}(\cdot|s)} [\tilde{Q}_{\pi_{tar}}(s, a) - \alpha \log \pi_{tar}(a|s) + \mu] \\ &= \tilde{V}_{\pi_{tar}}(s) + \mu. \end{aligned}$$

As $\pi_{tar}(\cdot|s)$ is contained in Π^s , with $\tilde{\pi}_g$ defined in Eq. (9), it is obvious that for all $s \in \mathcal{S}$,

$$\begin{aligned} & \mathbb{E}_{a \sim \tilde{\pi}_g(\cdot|s)} [\tilde{Q}_{\pi_{tar}}(s, a) - \alpha \log \tilde{\pi}_g(a|s)] \geq \\ & \mathbb{E}_{a \sim \pi_{tar}(\cdot|s)} [\tilde{Q}_{\pi_{tar}}(s, a) - \alpha \log \pi_{tar}(a|s)] = \tilde{V}_{\pi_{tar}}(s). \end{aligned} \tag{1}$$

Then for all $s_i \in \mathcal{S}$,

$$\begin{aligned} & V_{\pi_{tar}}(s_i) \leq \tilde{V}_{\pi_{tar}}(s_i) + \mu \\ &\leq \mathbb{E}_{a_i \sim \tilde{\pi}_g(\cdot|s_i)} [\tilde{Q}_{\pi_{tar}}(s_i, a_i) - \alpha \log \tilde{\pi}_g(a_i|s_i)] + \mu \\ &\leq \mathbb{E}_{a_i \sim \tilde{\pi}_g(\cdot|s_i)} [Q_{\pi_{tar}}(s_i, a_i) - \alpha \log \tilde{\pi}_g(a_i|s_i)] + 2\mu \\ &= \mathbb{E}_{a_i \sim \tilde{\pi}_g(a|s_i)} [r(s_i, a_i) - \alpha \log \tilde{\pi}_g(a_i|s_i) + \gamma V_{\pi_{tar}}(s_{i+1})] \\ &\quad + 2\mu \\ &\vdots \\ &\leq \mathbb{E}_{\tilde{\pi}_g} \left[\sum_{\tau=0}^{\infty} \gamma^\tau (r(s_{i+\tau}, a_{i+\tau}) - \alpha \log \tilde{\pi}_g(a_{i+\tau}|s_{i+\tau})) \right] \\ &\quad + 2 \sum_{\tau=0}^{\infty} \gamma^\tau \mu \\ &= V_{\tilde{\pi}_g}(s_i) + \frac{2\mu}{1-\gamma}. \end{aligned}$$

Proof of Theorem 2

Proof. According to the Pinsker's inequality [47], $D_{KL}(\pi_{tar}^{l+1}(\cdot|s) \|\tilde{\pi}_g^l(\cdot|s)) \geq \frac{1}{2 \ln 2} \|\pi_{tar}^{l+1}(\cdot|s) - \tilde{\pi}_g^l(\cdot|s)\|_1^2$, where $\|\cdot\|_1$ is the L1 norm. So we have that for all $s \in \mathcal{S}$, $\|\pi_{tar}^{l+1}(\cdot|s) - \tilde{\pi}_g^l(\cdot|s)\|_1 \leq \sqrt{2 \ln 2 \delta}$. According to the Performance Difference Lemma [48], we have that for all $s \in \mathcal{S}$:

$$V_{\tilde{\pi}_g^l}(s) - V_{\pi_{tar}^{l+1}}(s)$$

$$\begin{aligned}
&= \frac{1}{1-\gamma} \mathbb{E}_{s' \sim \rho_s^{\tilde{\pi}_g^l}}(s') [\mathbb{E}_{a \sim \tilde{\pi}_g^l(\cdot|s')} [Q_{\pi_{tar}^{l+1}}(s', a) - \alpha \log \tilde{\pi}_g^l(a|s)] \\
&\quad - \mathbb{E}_{a \sim \tilde{\pi}_{tar}^{l+1}(\cdot|s')} [Q_{\pi_{tar}^{l+1}}(s', a) - \alpha \log \tilde{\pi}_{tar}^{l+1}(a|s)]] \\
&\leq \frac{1}{1-\gamma} \max_{s' \in \mathcal{S}} [\mathbb{E}_{a \sim \tilde{\pi}_g^l(\cdot|s')} [Q_{\pi_{tar}^{l+1}}(s', a)] \\
&\quad - \mathbb{E}_{a \sim \tilde{\pi}_{tar}^{l+1}(\cdot|s')} [Q_{\pi_{tar}^{l+1}}(s', a)]] \\
&\quad + \frac{\alpha}{1-\gamma} \max_{s'' \in \mathcal{S}} |\mathcal{H}(\tilde{\pi}_g^l(\cdot|s'')) - \mathcal{H}(\pi_{tar}^l(\cdot|s''))| \\
&= \frac{1}{1-\gamma} \max_{s' \in \mathcal{S}} \int (\tilde{\pi}_g^l(\cdot|s) - \pi_{tar}^{l+1}(a|s)) Q_{\pi_{tar}^{l+1}}(s', a) da \\
&\quad + \frac{\alpha}{1-\gamma} \tilde{\mathcal{H}}_{max} \\
&\leq \frac{1}{1-\gamma} \max_{s' \in \mathcal{S}} \int |\tilde{\pi}_g^l(a|s) - \pi_{tar}^{l+1}(a|s)| \cdot |Q_{\pi_{tar}^{l+1}}(s', a)| da \\
&\quad + \frac{\alpha}{1-\gamma} \tilde{\mathcal{H}}_{max} \\
&\leq \frac{1}{1-\gamma} \max_{s' \in \mathcal{S}} \int |\tilde{\pi}_g^l(a|s) - \pi_{tar}^{l+1}(a|s)| \cdot \frac{\tilde{R}_{max} + \alpha \mathcal{H}_{max}^{l+1}}{1-\gamma} da \\
&\quad + \frac{\alpha}{1-\gamma} \tilde{\mathcal{H}}_{max} \\
&= \frac{\tilde{R}_{max} + \alpha \mathcal{H}_{max}^{l+1}}{(1-\gamma)^2} \max_{s' \in \mathcal{S}} \|\tilde{\pi}_g^l(\cdot|s) - \pi_{tar}^{l+1}(\cdot|s)\|_1 \\
&\quad + \frac{\alpha}{1-\gamma} \tilde{\mathcal{H}}_{max} \\
&\leq \frac{\sqrt{2 \ln 2\delta} (\tilde{R}_{max} + \alpha \mathcal{H}_{max}^{l+1}) + \alpha(1-\gamma) \tilde{\mathcal{H}}_{max}}{(1-\gamma)^2},
\end{aligned} \tag{2}$$

where $\rho_s^{\tilde{\pi}_g^l}(s') = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t p(s_t = s' | s_0 = s, \tilde{\pi}_g^l)$ is the normalized discounted state occupancy distribution. Note that

$$\begin{aligned}
&|Q_{\pi_{tar}^{l+1}}(s, a)| \\
&= |\mathbb{E}_{\pi_{tar}^{l+1}} [\sum_{i=0}^{\infty} \gamma^i (r(s_{\tau+i}, a_{\tau+i}) \\
&\quad - \alpha \log \pi_{tar}^{l+1}(\cdot|s)) | s_{\tau} = s, a_{\tau} = a]| \\
&\leq \mathbb{E}_{\pi} [\sum_{i=0}^{\infty} \gamma^i (\tilde{R}_{max} + \gamma \mathcal{H}_{max}^{l+1})]
\end{aligned} \tag{3}$$

$$= \frac{\tilde{R}_{max} + \alpha \mathcal{H}_{max}^{l+1}}{1 - \gamma}. \quad (4)$$

Eventually, we have

$$\begin{aligned} & V_{\pi_{tar}^{l+1}}(s) \\ & \geq V_{\tilde{\pi}_g^l}(s) - \frac{\sqrt{2 \ln 2 \delta} (\tilde{R}_{max} + \alpha \mathcal{H}_{max}^{l+1}) + \alpha(1 - \gamma) \tilde{\mathcal{H}}_{max}}{(1 - \gamma)^2} \\ & \geq V_{\pi_{tar}^l}(s) - \frac{\sqrt{2 \ln 2 \delta} (\tilde{R}_{max} + \alpha \mathcal{H}_{max}^{l+1})}{(1 - \gamma)^2} - \frac{2\mu + \alpha \tilde{\mathcal{H}}_{max}}{1 - \gamma}. \end{aligned} \quad (5)$$

References

- [1] Guberman, S.R., Greenfield, P.M.: Learning and transfer in everyday cognition. *Cognitive Development* **6**(3), 233–260 (1991)
- [2] Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., *et al.*: Mastering the game of go without human knowledge. *nature* **550**(7676), 354–359 (2017)
- [3] Vinyals, O., Babuschkin, I., Czarnecki, W.M., Mathieu, M., Dudzik, A., Chung, J., Choi, D.H., Powell, R., Ewalds, T., Georgiev, P., *et al.*: Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature* **575**(7782), 350–354 (2019)
- [4] Ceron, J.S.O., Castro, P.S.: Revisiting rainbow: Promoting more insightful and inclusive deep reinforcement learning research. In: *International Conference on Machine Learning*, pp. 1373–1383 (2021). PMLR
- [5] Fernández, F., Veloso, M.: Probabilistic policy reuse in a reinforcement learning agent. In: *Proceedings of the Fifth International Joint Conference on Autonomous Agents and Multiagent Systems*, pp. 720–727 (2006)
- [6] Barreto, A., Borsa, D., Quan, J., Schaul, T., Silver, D., Hessel, M., Mankowitz, D., Zidek, A., Munos, R.: Transfer in deep reinforcement learning using successor features and generalised policy improvement. In: *International Conference on Machine Learning*, pp. 501–510 (2018). PMLR
- [7] Li, S., Wang, R., Tang, M., Zhang, C.: Hierarchical reinforcement learning with advantage-based auxiliary rewards. *Advances in Neural Information Processing Systems* **32** (2019)
- [8] Yang, T., Hao, J., Meng, Z., Zhang, Z., Hu, Y., Chen, Y., Fan, C., Wang, W., Liu, W., Wang, Z., Peng, J.: Efficient deep reinforcement learning via adaptive policy transfer. In: *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, pp. 3094–3100 (2020)

- [9] Zhang, J., Li, S., Zhang, C.: Cup: Critic-guided policy reuse. In: *Advances in Neural Information Processing Systems* (2022)
- [10] Li, S., Gu, F., Zhu, G., Zhang, C.: Context-aware policy reuse. *arXiv preprint arXiv:1806.03793* (2018)
- [11] Kurenkov, A., Mandelkar, A., Martin-Martin, R., Savarese, S., Garg, A.: Ac-teach: A bayesian actor-critic method for policy learning with an ensemble of suboptimal teachers. In: *Conference on Robot Learning*, pp. 717–734 (2020). PMLR
- [12] Barreto, A., Dabney, W., Munos, R., Hunt, J.J., Schaul, T., Hasselt, H.P., Silver, D.: Successor features for transfer in reinforcement learning. *Advances in neural information processing systems* **30** (2017)
- [13] Cheng, C.-A., Kolobov, A., Agarwal, A.: Policy improvement via imitation of multiple oracles. *Advances in Neural Information Processing Systems* **33**, 5587–5598 (2020)
- [14] Pateria, S., Subagdja, B., Tan, A.-h., Quek, C.: Hierarchical reinforcement learning: A comprehensive survey. *ACM Computing Surveys (CSUR)* **54**(5), 1–35 (2021)
- [15] Lillicrap, T.P., Hunt, J.J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., Wierstra, D.: Continuous control with deep reinforcement learning. In: *ICLR (Poster)* (2016)
- [16] Fujimoto, S., Hoof, H., Meger, D.: Addressing function approximation error in actor-critic methods. In: *International Conference on Machine Learning*, pp. 1587–1596 (2018). PMLR
- [17] Haarnoja, T., Zhou, A., Hartikainen, K., Tucker, G., Ha, S., Tan, J., Kumar, V., Zhu, H., Gupta, A., Abbeel, P., et al.: Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905* (2018)
- [18] Yu, T., Quillen, D., He, Z., Julian, R., Hausman, K., Finn, C., Levine, S.: Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In: *Conference on Robot Learning*, pp. 1094–1100 (2020). PMLR
- [19] Zhu, Z., Lin, K., Zhou, J.: Transfer learning in deep reinforcement learning: A survey. *arXiv preprint arXiv:2009.07888* (2020)
- [20] Parisotto, E., Ba, J.L., Salakhutdinov, R.: Actor-mimic: Deep multitask and transfer reinforcement learning. *arXiv preprint arXiv:1511.06342* (2015)
- [21] Hou, Y., Ong, Y.-S., Feng, L., Zurada, J.M.: An evolutionary transfer reinforcement learning framework for multiagent systems. *IEEE Transactions on Evolutionary Computation* **21**(4), 601–615 (2017)

- [22] Larocche, R., Barlier, M.: Transfer reinforcement learning with shared dynamics. In: Thirty-First AAAI Conference on Artificial Intelligence (2017)
- [23] Lehnert, L., Littman, M.L.: Successor features combine elements of model-free and model-based reinforcement learning. *Journal of Machine Learning Research* **21**(196), 1–53 (2020)
- [24] Barekatin, M., Yonetani, R., Hamaya, M.: Multipolar: multi-source policy aggregation for transfer reinforcement learning between diverse environmental dynamics. In: Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence, pp. 3108–3116 (2020)
- [25] Li, S., Zhang, C.: An optimal online method of selecting source policies for reinforcement learning. In: Proceedings of the AAAI Conference on Artificial Intelligence, vol. 32 (2018)
- [26] Gimelfarb, M., Sanner, S., Lee, C.-G.: Contextual policy transfer in reinforcement learning domains via deep mixtures-of-experts. In: Uncertainty in Artificial Intelligence, pp. 1787–1797 (2021). PMLR
- [27] Yang, X., Ji, Z., Wu, J., Lai, Y.-K., Wei, C., Liu, G., Setchi, R.: Hierarchical reinforcement learning with universal policies for multistep robotic manipulation. *IEEE Transactions on Neural Networks and Learning Systems* **33**(9), 4727–4741 (2021)
- [28] Rusu, A.A., Rabinowitz, N.C., Desjardins, G., Soyer, H., Kirkpatrick, J., Kavukcuoglu, K., Pascanu, R., Hadsell, R.: Progressive neural networks. arXiv preprint arXiv:1606.04671 (2016)
- [29] Berseth, G., Xie, C., Cernek, P., Panne, M.: Progressive reinforcement learning with distillation for multi-skilled motion control. arXiv preprint arXiv:1802.04765 (2018)
- [30] Schwarz, J., Czarnecki, W., Luketina, J., Grabska-Barwinska, A., Teh, Y.W., Pascanu, R., Hadsell, R.: Progress & compress: A scalable framework for continual learning. In: International Conference on Machine Learning, pp. 4528–4537 (2018). PMLR
- [31] Mallya, A., Lazebnik, S.: Packnet: Adding multiple tasks to a single network by iterative pruning. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 7765–7773 (2018)
- [32] Teh, Y., Bapst, V., Czarnecki, W.M., Quan, J., Kirkpatrick, J., Hadsell, R., Heess, N., Pascanu, R.: Distral: Robust multitask reinforcement learning. *Advances in neural information processing systems* **30** (2017)
- [33] Haarnoja, T., Zhou, A., Abbeel, P., Levine, S.: Soft actor-critic: Off-policy

- maximum entropy deep reinforcement learning with a stochastic actor. In: International Conference on Machine Learning, pp. 1861–1870 (2018). PMLR
- [34] Lan, Q., Pan, Y., Fyshe, A., White, M.: Maxmin q-learning: Controlling the estimation bias of q-learning. arXiv preprint arXiv:2002.06487 (2020)
- [35] Kuznetsov, A., Shvechikov, P., Grishin, A., Vetrov, D.: Controlling overestimation bias with truncated mixture of continuous distributional quantile critics. In: International Conference on Machine Learning, pp. 5556–5566 (2020). PMLR
- [36] Zhang, S., Sutton, R.S.: A deeper look at experience replay. arXiv preprint arXiv:1712.01275 (2017)
- [37] Fedus, W., Ramachandran, P., Agarwal, R., Bengio, Y., Laroche, H., Rowland, M., Dabney, W.: Revisiting fundamentals of experience replay. In: International Conference on Machine Learning, pp. 3061–3071 (2020). PMLR
- [38] Khetarpal, K., Riemer, M., Rish, I., Precup, D.: Towards continual reinforcement learning: A review and perspectives. arXiv preprint arXiv:2012.13490 (2020)
- [39] Wolczyk, M., Zajkac, M., Pascanu, R., Kucinski, L., Milos, P.: Continual world: A robotic benchmark for continual reinforcement learning. *Advances in Neural Information Processing Systems* **34**, 28496–28510 (2021)
- [40] Yang, R., Xu, H., Wu, Y., Wang, X.: Multi-task reinforcement learning with soft modularization. *Advances in Neural Information Processing Systems* **33**, 4767–4777 (2020)
- [41] Sodhani, S., Zhang, A., Pineau, J.: Multi-task reinforcement learning with context-based representations. In: International Conference on Machine Learning, pp. 9767–9779 (2021). PMLR
- [42] Wan, M., Gangwani, T., Peng, J.: Mutual information based knowledge transfer under state-action dimension mismatch. arXiv preprint arXiv:2006.07041 (2020)
- [43] Zhang, Q., Xiao, T., Efros, A.A., Pinto, L., Wang, X.: Learning cross-domain correspondence for control with dynamics cycle-consistency. arXiv preprint arXiv:2012.09811 (2020)
- [44] Heng, Y., Yang, T., ZHENG, Y., Jianye, H., Taylor, M.E.: Cross-domain adaptive transfer reinforcement learning based on state-action correspondence. In: The 38th Conference on Uncertainty in Artificial Intelligence (2022)
- [45] Pol, E., Worrall, D., Hoof, H., Oliehoek, F., Welling, M.: Mdp homomorphic networks: Group symmetries in reinforcement learning. *Advances in Neural Information Processing Systems* **33**, 4199–4210 (2020)
- [46] Pol, E., Kipf, T., Oliehoek, F.A., Welling, M.: Plannable approximations to mdp

homomorphisms: Equivariance under actions. arXiv preprint arXiv:2002.11963 (2020)

- [47] Fedotov, A.A., Harremoës, P., Topsøe, F.: Refinements of pinsker's inequality. *IEEE Transactions on Information Theory* **49**(6), 1491–1498 (2003)
- [48] Kakade, S., Langford, J.: Approximately optimal approximate reinforcement learning. In: *In Proc. 19th International Conference on Machine Learning* (2002). Citeseer

