# Learning Action Translator for Meta Reinforcement Learning on Sparse-Reward Tasks

**Yijie Guo[1], Qiucheng Wu[1], Honglak Lee[1,2]**

[1]University of Michigan, [2]LG AI Research
guoyijie@umich.edu, wuqiuche@umich.edu, honglak@eecs.umich.edu

## Abstract

Meta reinforcement learning (meta-RL) aims to learn a policy solving a set of training tasks simultaneously and quickly adapting to new tasks. It requires massive amounts of data drawn from training tasks to infer the common structure shared among tasks. Without heavy reward engineering, the sparse rewards in long-horizon tasks exacerbate the problem of sample efficiency in meta-RL. Another challenge in meta-RL is the discrepancy of difficulty level among tasks, which might cause one easy task dominating learning of the shared policy and thus preclude policy adaptation to new tasks. This work introduces a novel objective function to learn an action translator among training tasks. We theoretically verify that the value of the transferred policy with the action translator can be close to the value of the source policy and our objective function (approximately) upper bounds the value difference. We propose to combine the action translator with context-based meta-RL algorithms for better data collection and more efficient exploration during meta-training. Our approach empirically improves the sample efficiency and performance of meta-RL algorithms on sparse-reward tasks.

## 1 Introduction

Deep reinforcement learning (DRL) methods achieved remarkable success in solving complex tasks (Mnih et al. 2015; Silver et al. 2016; Schulman et al. 2017). While conventional DRL methods learn an individual policy for each task, meta reinforcement learning (meta-RL) algorithms (Finn, Abbeel, and Levine 2017; Duan et al. 2016; Mishra et al. 2017) learn the shared structure across a distribution of tasks so that the agent can quickly adapt to unseen related tasks in the test phase. Unlike most of the existing meta-RL approaches working on tasks with dense rewards, we instead focus on the sparse-reward training tasks, which are more common in real-world scenarios without access to carefully designed reward functions in the environments. Recent works in meta-RL propose off-policy algorithms (Rakelly et al. 2019; Fakoor et al. 2019) and model-based algorithms (Nagabandi, Finn, and Levine 2018; Nagabandi et al. 2018) to improve the sample efficiency in meta-training procedures. However, it remains challenging to efficiently solve multiple tasks that require reasoning over long
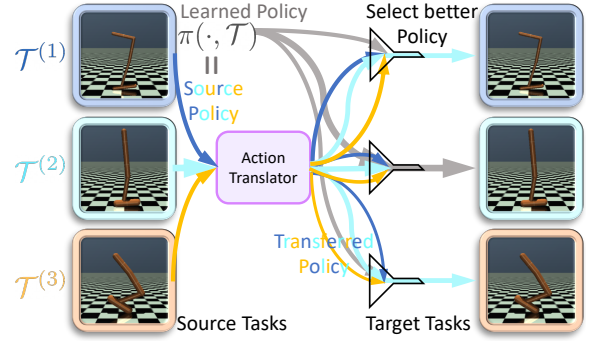
Figure 1: Illustration of our policy transfer. Size of arrows represents avg. episode reward of learned or transferred policy on target tasks. Different colors indicate different tasks.

horizons with sparse rewards. In these tasks, the scarcity of positive rewards exacerbates the issue of sample efficiency, which plagues meta-RL algorithms and makes exploration difficult due to a lack of guidance signals.

Intuitively, we hope that solving one task facilitates learning of other related tasks since the training tasks share a common structure. However, it is often not the case in practice (Rusu et al. 2015; Parisotto, Ba, and Salakhutdinov 2015). Previous works (Teh et al. 2017; Yu et al. 2020a) point out that detrimental gradient interference might cause an imbalance in policy learning on multiple tasks. Policy distillation (Teh et al. 2017) and gradient projection (Yu et al. 2020a) are developed in meta-RL algorithms to alleviate this issue. However, this issue might become more severe in the sparse-reward setting because it is hard to explore each task to obtain meaningful gradient signals for policy updates. Good performance in one task does not automatically help exploration on the other tasks since the agent lacks positive rewards on the other tasks to learn from.

In this work, we aim to fully exploit the highly-rewarding transitions occasionally discovered by the agent in the exploration. The good experiences in one task should not only improve the policy on this task but also benefit the policy on other tasks to drive deeper exploration. Specifically, once the agent learns from the successful trajectories in one training task, we transfer the good policy in this task to other tasks to

get more positive rewards on other training tasks. In Fig. 1, if the learned policy $\pi$ performs better on task $\mathcal{T}^{(2)}$ than other tasks, then our goal is to transfer the good policy $\pi(\cdot, \mathcal{T}^{(2)})$ to other tasks $\mathcal{T}^{(1)}$ and $\mathcal{T}^{(3)}$. To enable such transfer, we propose to learn an action translator among multiple training tasks. The objective function forces the translated action to behave on the target task similarly to the source action on the source task. We consider the policy transfer for any pair of source and target tasks in the training task distribution (see the colored arrows in Fig. 1). The agent executes actions following the transferred policy if the transferred policy attains higher rewards than the learned policy on the target task in recent episodes. This approach enables the agent to leverage relevant data from multiple training tasks, encourages the learned policy to perform similarly well on multiple training tasks, and thus leads to better performance when applying the well-trained policy to test tasks.

We summarize the contributions: (1) We introduce a novel objective function to transfer any policy from a source Markov Decision Process (MDP) to a target MDP. We prove a theoretical guarantee that the transferred policy can achieve the expected return on the target MDP close to the source policy on the source MDP. The difference in expected returns is (approximately) upper bounded by our loss function with a constant multiplicative factor. (2) We develop an off-policy RL algorithm called **M**eta-RL with **C**ontext-conditioned **A**ction **T**ranslator (MCAT), applying a policy transfer mechanism in meta-RL to help exploration across multiple sparse-rewards tasks. (3) We empirically demonstrate the effectiveness of MCAT on a variety of simulated control tasks with the MuJoCo physics engine (Todorov, Erez, and Tassa 2012), showing that policy transfer improves the performance of context-based meta-RL algorithms.

## 2 Related Work

**Context-based Meta-RL** Meta reinforcement learning has been extensively studied in the literature (Finn, Abbeel, and Levine 2017; Stadie et al. 2018; Sung et al. 2017; Xu, van Hasselt, and Silver 2018) with many works developing the context-based approaches (Rakelly et al. 2019; Ren et al. 2020; Liu et al. 2020). Duan et al. (2016); Wang et al. (2016); Fakoor et al. (2019) employ recurrent neural networks to encode context transitions and formulate the policy conditioning on the context variables. The objective function of maximizing expected return trains the context encoder and policy jointly. Rakelly et al. (2019) leverage a permutation-invariant encoder to aggregate experiences as probabilistic context variables and optimizes it with variational inference. The posterior sampling is beneficial for exploration on sparse-reward tasks in the adaptation phase, but there is access to dense rewards during training phase. Li, Pinto, and Abbeel (2020) considers a task-family of reward functions. Lee et al. (2020); Seo et al. (2020) trains the context encoder with forward dynamics prediction. These model-based meta-RL algorithms assume the reward function is accessible for planning. In the sparse-reward setting without ground-truth reward functions, they may struggle to discover non-zero rewards and accurately estimating the reward for model-based planning may be problematic as well.

**Policy Transfer in RL** Policy transfer studies the knowledge transfer in target tasks given a set of source tasks and their expert policies. Policy distillation (Rusu et al. 2015; Yin and Pan 2017; Parisotto, Ba, and Salakhutdinov 2015) minimize the divergence of action distributions between the source policy and the learned policy on the target task. Along this line of works, Teh et al. (2017) create a centroid policy in multi-task reinforcement learning and distills the knowledge from the task-specific policies to this centroid policy. Alternatively, inter-task mapping between the source and target tasks (Zhu, Lin, and Zhou 2020) can assist the policy transfer. Most of these works (Gupta et al. 2017; Konidaris and Barto 2006; Ammar and Taylor 2011) assume existence of correspondence over the state space and learn the state mapping between tasks. Recent work (Zhang et al. 2020c) learns the state correspondence and action correspondence with dynamic cycle-consistency loss. Our method differs from this approach, in that we enable action translation among multiple tasks with a simpler objective function. Importantly, our approach is novel to utilize the policy transfer for any pair of source and target tasks in meta-RL.

**Bisimulation for States in MDPs** Recent works on state representation learning (Ferns, Panangaden, and Precup 2004; Zhang et al. 2020a; Agarwal et al. 2021) investigate the bismilarity metrics for states on multiple MDPs and consider how to learn a representation for states leading to almost identical behaviors under the same action in diverse MDPs. In multi-task reinforcement learning and meta reinforcement learning problems, Zhang et al. (2020a,b) derives transfer and generalization bounds based on the task and state similarity. We also bound the value of policy transfer across tasks but our approach is to establish action equivalence instead of state equivalence.

## 3 Method

In this section, we first describe our approach to learn a context encoder capturing the task features and learn a forward dynamics model predicting next state distribution given the task context (Sec. 3.2). Then we introduce an objective function to train an action translator so that the translated action on the target task behaves equivalently to the source action on the source task. The action translator can be conditioned on the task contexts and thus it can transfer a good policy from any arbitrary source task to any other target task in the training set (Sec. 3.3). Finally, we propose to combine the action translator with a context-based meta-RL algorithm to transfer the good policy from any one task to the others. During meta-training, this policy transfer approach helps exploit the good experiences encountered on any one task and benefits the data collection and further policy optimization on other sparse-reward tasks (Sec. 3.4). Fig. 2 provides an overview of our approach MCAT.

### 3.1 Problem Formulation

Following meta-RL formulation in previous work (Duan et al. 2016; Mishra et al. 2017; Rakelly et al. 2019), we assume a distribution of tasks $p(\mathcal{T})$ and each task is a Markov decision process (MDP) defined as a tuple $(\mathcal{S}, \mathcal{A}, p, r, \gamma, \rho_0)$
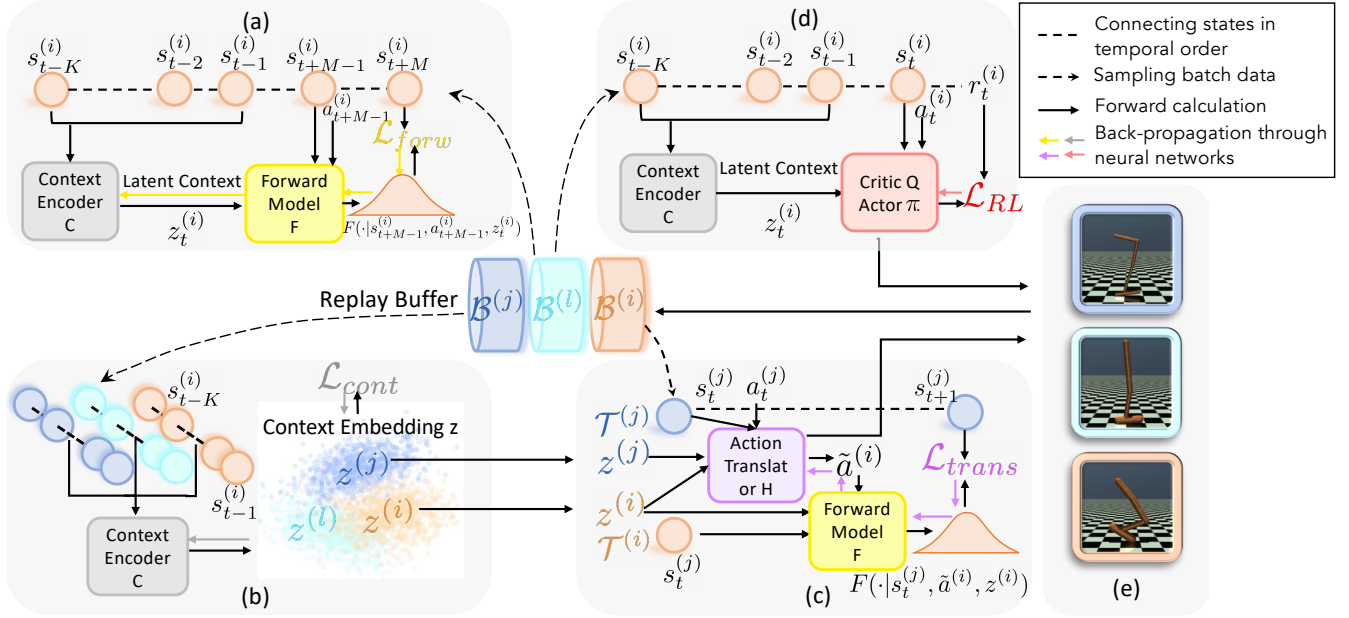
Figure 2: Overview of MCAT. (a) We use forward dynamics prediction loss to train the context encoder $C$ and forward model $F$. (b) We regularize the context encoder $C$ with the contrastive loss, so context vectors of transition segments from the same task cluster together. (c) With fixed $C$ and $F$, we learn the action translator $H$ for any pair of source task $\mathcal{T}^{(j)}$ and target task $\mathcal{T}^{(i)}$. The action translator aims to generate action $\tilde{a}^{(i)}$ on the target task leading to the same next state $s_{t+1}^{(j)}$ as the source action $a_t^{(j)}$ on the source task. (d) With fixed $C$, we learn the critic $Q$ and actor $\pi$ conditioning on the context feature. (e) If the agent is interacting with the environment on task $\mathcal{T}^{(i)}$, we compare learned policy $\pi(s, z^{(i)})$ and transferred policy $H(s, \pi(s, z^{(j)}), z^{(j)}, z^{(i)})$, which transfers a good policy $\pi(s, z^{(j)})$ on source task $\mathcal{T}^{(j)}$ to target task $\mathcal{T}^{(i)}$. We select actions according to the policy with higher average episode rewards in the recent episodes. Transition data are pushed into the buffer. We remark that the components $C, F, H, Q, \pi$ are trained alternatively not jointly and this fact facilitates the learning process.

with state space $\mathcal{S}$, action space $\mathcal{A}$, transition function $p(s'|s,a)$, reward function $r(s,a,s')$, discounting factor $\gamma$, and initial state distribution $\rho_0$. We can alternatively define the reward function as $r(s,a) = \sum_{s' \in \mathcal{S}} p(s'|s,a)r(s,a,s')$. In context-based meta-RL algorithms, we learn a policy $\pi(\cdot|s_t^{(i)}, z_t^{(i)})$ shared for any task $\mathcal{T}^{(i)} \sim p(\mathcal{T})$, where $t$ denotes the timestep in an episode, $i$ denotes the index of a task, the context variable $z_t^{(i)} \in \mathcal{Z}$ captures contextual information from history transitions on the task MDP and $\mathcal{Z}$ is the space of context vectors. The shared policy is optimized to maximize its value $V^\pi(\mathcal{T}^{(i)}) = \mathbb{E}_{\rho_0^{(i)}, \pi, p^{(i)}}[\sum_{t=0}^\infty \gamma^t r_t^{(i)}]$ on each training task $\mathcal{T}^{(i)}$. Following prior works in meta-RL (Yu et al. 2017; Nagabandi et al. 2018; Nagabandi, Finn, and Levine 2018; Zhou, Pinto, and Gupta 2019; Lee et al. 2020), we study tasks with the same state space, action space, reward function but varying dynamics functions. Importantly, we focus on more challenging setting of sparse rewards. Our goal is to learn a shared policy robust to the dynamic changes and generalizable to unseen tasks.

### 3.2 Learning Context & Forward Model

In order to capture the knowledge about any task $\mathcal{T}^{(i)}$, we leverage a context encoder $C : \mathcal{S}^K \times \mathcal{A}^K \to \mathcal{Z}$, where $K$ is the number of past steps used to infer the context. Related ideas have been explored by (Rakelly et al. 2019; Zhou, Pinto, and Gupta 2019; Lee et al. 2020). In Fig. 2a,

given $K$ past transitions $(s_{t-K}^{(i)}, a_{t-K}^{(i)}, \cdots, s_{t-1}^{(i)}, a_{t-1}^{(i)})$, context encoder $C$ produces the latent context $z_t^{(i)} = C(s_{t-K}^{(i)}, a_{t-K}^{(i)}, \cdots, s_{t-2}^{(i)}, a_{t-2}^{(i)}, s_{t-1}^{(i)}, a_{t-1}^{(i)})$. We train the context encoder $C$ and forward dynamics $F$ with an objective function to predict the forward dynamics in future transitions $s_{t+m}^{(i)}$ ($1 \leq m \leq M$) within $M$ future steps. The state prediction in multiple future steps drives latent context embeddings $z_t^{(i)}$ to be temporally consistent. The learned context encoder tends to capture dynamics-specific, contextual information (e.g. environment physics parameters). Formally, we minimize the negative log-likelihood of observing the future states under dynamics prediction.

$$\mathcal{L}_{forw} = -\sum_{m=1}^M \log F(s_{t+m}^{(i)}|s_{t+m-1}^{(i)}, a_{t+m-1}^{(i)}, z_t^{(i)}). \quad (1)$$

Additionally, given trajectory segments from the same task, we require their context embeddings to be similar, whereas the contexts of history transitions from different tasks should be distinct (Fig. 2b). We propose a contrastive loss (Hadsell, Chopra, and LeCun 2006) to constrain embeddings within a small distance for positive pairs (i.e. samples from the same task) and push embeddings apart with a distance greater than a margin value $m$ for negative pairs (i.e. samples from different tasks). $z_{t_1}^{(i)}, z_{t_2}^{(j)}$ denote context embeddings of two trajectory samples from $\mathcal{T}^{(i)}, \mathcal{T}^{(j)}$. The

contrastive loss function is defined as:

$$\mathcal{L}_{cont} = 1_{i=j}\|z_{t_1}^{(i)} - z_{t_2}^{(j)}\|^2 + 1_{i \neq j}\max(0, m - \|z_{t_1}^{(i)} - z_{t_2}^{(j)}\|) \tag{2}$$

where 1 is indicator function. During meta-training, recent transitions on each task $\mathcal{T}^{(i)}$ are stored in a buffer $\mathcal{B}^{(i)}$ for off-policy learning. We randomly sample a fairly large batch of trajectory segments from $\mathcal{B}^{(i)}$, and average their context embeddings to output task feature $z^{(i)}$. $z^{(i)}$ is representative for embeddings on task $\mathcal{T}^{(i)}$ and distinctive from features $z^{(l)}$ and $z^{(j)}$ for other tasks. We note the learned embedding maintains the similarity across tasks. $z^{(i)}$ is closer to $z^{(l)}$ than to $z^{(j)}$ if task $\mathcal{T}^{(i)}$ is more akin to $\mathcal{T}^{(l)}$. We utilize task features for action translation across multiple tasks. Appendix D.5 visualizes context embeddings to study $\mathcal{L}_{cont}$.

### 3.3 Learning Action Translator

Suppose that transition data $s_t^{(j)}, a_t^{(j)}, s_{t+1}^{(j)}$ behave well on task $\mathcal{T}^{(j)}$. We aim to learn an action translator $H : \mathcal{S} \times \mathcal{A} \times \mathcal{Z} \times \mathcal{Z} \rightarrow \mathcal{A}$. $\tilde{a}^{(i)} = H(s_t^{(j)}, a_t^{(j)}, z^{(j)}, z^{(i)})$ translates the proper action $a_t^{(j)}$ from source task $\mathcal{T}^{(j)}$ to target task $\mathcal{T}^{(i)}$. In Fig. 2c, if we start from the same state $s_t^{(j)}$ on both source and target tasks, the translated action $\tilde{a}^{(i)}$ on target task should behave equivalently to the source action $a_t^{(j)}$ on the source task. Thus, the next state $s_{t+1}^{(i)} \sim p^{(i)}(s_t^{(j)}, \tilde{a}^{(i)})$ produced from the transferred action $\tilde{a}^{(i)}$ on the target task should be close to the real next state $s_{t+1}^{(j)}$ gathered on the source task. The objective function of training the action translator $H$ is to maximize the probability of getting next state $s_{t+1}^{(j)}$ under the next state distribution $s_{t+1}^{(i)} \sim p^{(i)}(s_t^{(j)}, \tilde{a}^{(i)})$ on the target task. Because the transition function $p^{(i)}(s_t^{(j)}, \tilde{a}^{(i)})$ is unavailable and might be not differentiable, we use the forward dynamics model $F(\cdot|s_t^{(j)}, \tilde{a}^{(i)}, z^{(i)})$ to approximate the transition function. We formulate objective function for action translator $H$ as:

$$\mathcal{L}_{trans} = -\log F(s_{t+1}^{(j)}|s_t^{(j)}, \tilde{a}^{(i)}, z^{(i)}) \tag{3}$$

where $\tilde{a}^{(i)} = H(s_t^{(j)}, a_t^{(j)}, z^{(j)}, z^{(i)})$. We assume to start from the same initial state, the action translator is to find the action on the target task so as to reach the same next state as the source action on the source task. This intuition to learn the action translator is analogous to learn inverse dynamic model across two tasks.

With a well-trained action translator conditioning on task features $z^{(j)}$ and $z^{(i)}$, we transfer the good deterministic policy $\pi(s, z^{(j)})$ from any source task $\mathcal{T}^{(j)}$ to any target task $\mathcal{T}^{(i)}$. When encountering a state $s^{(i)}$ on $\mathcal{T}^{(i)}$, we query a good action $a^{(j)} = \pi(s^{(i)}, z^{(j)})$ which will lead to a satisfactory next state with high return on the source task. Then $H$ translates this good action $a^{(j)}$ on the source task to action $\tilde{a}^{(i)} = H(s^{(i)}, a^{(j)}, z^{(j)}, z^{(i)})$ on the target task. Executing the translated action $\tilde{a}^{(i)}$ moves the agent to a next state on the target task similarly to the good action on the source task. Therefore, transferred policy $H(s^{(i)}, \pi(s^{(i)}, z^{(j)}), z^{(i)}, z^{(j)})$ can behave similarly to source policy $\pi(s, z^{(j)})$. Sec. 5.1 demonstrates the performance of transferred policy in a variety of environments. Our policy transfer mechanism is related to the action correspondence discussed in (Zhang et al.

2020c). We extend their policy transfer approach across two domains to multiple domains(tasks) and theoretically validate learning of action translator in Sec. 4.

### 3.4 Combining with Context-based Meta-RL

MCAT follows standard off-policy meta-RL algorithms to learn a deterministic policy $\pi(s_t, z_t^{(i)})$ and a value function $Q(s_t, a_t, z_t^{(i)})$, conditioning on the latent task context variable $z_t^{(i)}$. In the meta-training process, using data sampled from $\mathcal{B}$, we train the context model $C$ and dynamics model $F$ with $\mathcal{L}_{forw}$ and $\mathcal{L}_{cont}$ to accurately predict the next state (Fig. 2a 2b). With the fixed context encoder $C$ and dynamics model $F$, the action translator $H$ is optimized to minimize $\mathcal{L}_{trans}$ (Fig. 2c). Then, with the fixed $C$, we train the context-conditioned policy $\pi$ and value function $Q$ according to $\mathcal{L}_{RL}$ (Fig. 2d). In experiments, we use the objective function $\mathcal{L}_{RL}$ from TD3 algorithm (Fujimoto, Hoof, and Meger 2018). See pseudo-code of MCAT in Appendix B.

On sparse-reward tasks where exploration is challenging, the agent might luckily find transitions with high rewards on one task $\mathcal{T}^{(j)}$. Thus, the policy learning on this task might be easier than other tasks. If the learned policy $\pi$ performs better on one task $\mathcal{T}^{(j)}$ than another task $\mathcal{T}^{(i)}$, we consider the policy transferred from $\mathcal{T}^{(j)}$ to $\mathcal{T}^{(i)}$. At a state $s^{(i)}$, we employ the action translator to get a potentially good action $H(s^{(i)}, \pi(s^{(i)}, z^{(j)}), z^{(j)}, z^{(i)})$ on target task $\mathcal{T}^{(i)}$. As illustrated in Fig. 2e and Fig. 1, in the recent episodes, if the transferred policy earns higher scores than the learned policy $\pi(s^{(i)}, z^{(i)})$ on the target task $\mathcal{T}^{(i)}$, we follow the translated actions on $\mathcal{T}^{(i)}$ to gather transition data in the current episode. These data with better returns are pushed into the replay buffer $\mathcal{B}^{(i)}$ and produce more positive signals for policy learning in the sparse-reward setting. These transition samples help improve $\pi$ on $\mathcal{T}^{(i)}$ after policy update with off-policy RL algorithms. As described in Sec. 3.3, our action translator $H$ allows policy transfer across any pair of tasks. Therefore, with the policy transfer mechanism, the learned policy on each task might benefit from good experiences and policies on any other tasks.

## 4 Theoretical Analysis

In this section, we theoretically support our objective function (Equation 3) to learn the action translator. Given $s$ on two MDPs with the same state and action space, we define that action $a^{(i)}$ on $\mathcal{T}^{(i)}$ is equivalent to action $a^{(j)}$ on $\mathcal{T}^{(j)}$ if the actions yielding exactly the same next state distribution and reward, i.e. $p^{(i)}(\cdot|s, a^{(i)}) = p^{(j)}(\cdot|s, a^{(j)})$ and $r^{(i)}(s, a^{(i)}) = r^{(j)}(s, a^{(j)})$. Ideally, the equivalent action always exists on the target MDP $\mathcal{T}^{(i)}$ for any state-action pair on the source MDP $\mathcal{T}^{(j)}$ and there exists an action translator function $H : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{A}$ to identify the exact equivalent action. Starting from state $s$, the translated action $\tilde{a} = H(s, a)$ on the task $\mathcal{T}^{(i)}$ generates reward and next state distribution the same as action $a$ on the task $\mathcal{T}^{(j)}$ (i.e. $\tilde{a}B_s a$). Then any deterministic policy $\pi^{(j)}$ on the source task $\mathcal{T}^{(j)}$ can be perfectly transferred to the target task $\mathcal{T}^{(i)}$ with $\pi^{(i)}(s) = H(s, \pi^{(j)}(s))$. The value of the policy $\pi^{(j)}$ on the source task $\mathcal{T}^{(j)}$ is equal to the value of transferred policy

$\pi^{(i)}$ on the target task $\mathcal{T}^{(i)}$.

Without the assumption of existence of a perfect correspondence for each action, given any two deterministic policies $\pi^{(j)}$ on $\mathcal{T}^{(j)}$ and $\pi^{(i)}$ on $\mathcal{T}^{(i)}$, we prove that the difference in the policy value is upper bounded by a scalar $\frac{d}{1-\gamma}$ depending on L1-distance between reward functions $|r^{(i)}(s, \pi^{(i)}(s)) - r^{(j)}(s, \pi^{(j)}(s))|$ and total-variation distance between next state distributions $D_{TV}(p^{(i)}(\cdot|s, \pi^{(i)}(s)), p^{(j)}(\cdot|s, \pi^{(j)}(s)))$. Detailed theorem (Theorem 1) and proof are in Appendix A.

For a special case where reward function $r(s, a, s')$ only depends on the current state $s$ and next state $s'$, the upper bound of policy value difference is only related to the distance in next state distributions.

**Proposition 1.** *Let* $\mathcal{T}^{(i)} = \{\mathcal{S}, \mathcal{A}, p^{(i)}, r^{(i)}, \gamma, \rho_0\}$ *and* $\mathcal{T}^{(j)} = \{\mathcal{S}, \mathcal{A}, p^{(j)}, r^{(j)}, \gamma, \rho_0\}$ *be two MDPs sampled from the distribution of tasks* $p(\mathcal{T})$. $\pi^{(i)}$, $\pi^{(j)}$ *is the deterministic policy on* $\mathcal{T}^{(i)}$, $\mathcal{T}^{(j)}$. *Assume the reward function only depends on the state and next state* $r^{(i)}(s, a^{(i)}, s') = r^{(j)}(s, a^{(j)}, s') = r(s, s')$. *Let* $d = \sup_{s \in \mathcal{S}} 2M D_{TV}(p^{(j)}(\cdot|s, \pi^{(j)}(s)), p^{(i)}(\cdot|s, \pi^{(i)}(s)))$ *and* $M = \sup_{s \in \mathcal{S}, s' \in \mathcal{S}} |r(s, s') + \gamma V^{\pi^{(j)}}(s, \mathcal{T}^{(j)})|$. $\forall s \in \mathcal{S}$, *we have*

$$\left| V^{\pi^{(i)}}(s, \mathcal{T}^{(i)}) - V^{\pi^{(j)}}(s, \mathcal{T}^{(j)}) \right| \leq \frac{d}{1-\gamma} \quad (4)$$

According to Proposition 1, if we can optimize the action translator $H$ to minimize $d$ for policy $\pi^{(j)}$ and $\pi^{(i)}(s) = H(s, \pi^{(j)}(s))$, the value of the transferred policy $\pi^{(i)}$ on the target task can be close to the value of source policy $\pi^{(j)}$. In many real-world scenarios, especially sparse-reward tasks, the reward heavily depends on the state and next state instead of action. For example, robots running forward receive rewards according to their velocity (i.e. the location difference between the current and next state within one step); robot arms manipulating various objects earn positive rewards only when they are in the target positions. Thus, our approach focuses on the cases with reward function approximately as $r(s, s')$ under the assumption of Proposition 1. For any state $s \in \mathcal{S}$, we minimize the total-variation distance between two next state distributions $D_{TV}(p^{(j)}(\cdot|s_t, \pi^{(j)}(s_t)), p^{(i)}(\cdot|s_t, \pi^{(i)}(s_t)))$ on source and target MDPs. Besides, we discuss the policy transfer for tasks with a general reward function in Appendix C.3.

There is no closed-form solution of $D_{TV}$ and $D_{TV}$ is related with Kullback–Leibler (KL) divergence $D_{KL}$ by the inequality $D_{TV}(p\|q)^2 \leq D_{KL}(p\|q)$ Thus, we instead consider minimizing $D_{KL}$ between two next state distributions. $D_{KL}(p^{(j)}\|p^{(i)})$ is $-\sum_{s'} p^{(j)}(s') \log p^{(i)}(s') + \sum_{s'} p^{(j)}(s') \log p^{(j)}(s')$. The second term does not involve $H$ and thus can be viewed as a constant term when optimizing $H$. We focus on minimizing the first term $-\sum_{s'} p^{(j)}(s') \log p^{(i)}(s')$. $F$ is a forward model approximating $p^{(i)}(s')$. We sample transitions $s, \pi^{(j)}(s), s'$ from the source task. $s'$ follows the distribution $p^{(j)}(s')$. Thus, minimizing the negative log-likelihood of observing the next state $L_{trans} = -\log F(s'|s, \pi^{(i)}(s))$ is to approximately minimize $D_{KL}$. Experiments in Sec. 5.1 suggest that this objective function works well for policy transfer across

two MDPs. Sec. 3.3 explains the motivation behind $\mathcal{L}_{trans}$ (Equation 3) to learn an action translator among multiple MDPs instead of only two MDPs.

# 5 Experiment

We design and conduct experiments to answer the following questions: (1) Does the transferred policy perform well on the target task (Tab. 1, Fig. 4)? (2) Can we transfer the good policy for any pair of source and target tasks (Fig. 5)? (3) Does policy transfer improve context-based Meta-RL algorithms (Fig. 3, Tab. 2, Tab. 3)? (4) Is the policy transfer more beneficial when the training tasks have sparser rewards (Tab. 4)? Experimental details can be found in Appendix C.

## 5.1 Policy Transfer with Fixed Dataset

We test our proposed action translator with fixed datasets of transitions aggregated from pairs of source and target tasks. On MuJoCo environments HalfCheetah and Ant, we create tasks with varying dynamics as in (Zhou, Pinto, and Gupta 2019; Lee et al. 2020; Zhang et al. 2020c). We keep default physics parameters in source tasks and modify them to yield noticeable changes in the dynamics for target tasks. On HalfCheetah, the tasks differ in the armature. On Ant, we set different legs crippled. A well-performing policy is pretrained on the source task with TD3 algorithm (Fujimoto, Hoof, and Meger 2018) and dense rewards. We then gather training data with mediocre policies on the source and target tasks. We also include object manipulation tasks on Meta-World benchmark (Yu et al. 2020b). Operating objects with varied physics properties requires the agent to handle different dynamics. The knowledge in grasping and pushing a cylinder might be transferrable to tasks of moving a coffee mug or a cube. The agent gets a reward of 1.0 if the object is in the goal location. Otherwise, the reward is 0. We use the manually-designed good policy as the source policy and collect transition data by adding noise to the action drawn from the good policy.

| Setting | Source policy | Transferred policy (Zhang et al. 2020c) | Transferred policy (Ours) |
|---|---|---|---|
| HalfCheetah | 2355.0 | **3017.1**($\pm$44.2) | 2937.2($\pm$9.5) |
| Ant | 55.8 | 97.2($\pm$2.5) | **208.1**($\pm$8.2) |
| Cylinder-Mug | 0.0 | 308.1($\pm$75.3) | **395.6**($\pm$19.4) |
| Cylinder-Cube | 0.0 | 262.4($\pm$48.1) | **446.1**($\pm$1.1) |

Table 1: Mean ($\pm$ standard error) of episode rewards over 3 runs, comparing source and transferred policy on target task.

As presented in Tab. 1, directly applying a good source policy on the target task performs poorly. We learn dynamics model $F$ on target task with $\mathcal{L}_{forw}$ and action translator $H$ with $\mathcal{L}_{trans}$. From a single source task to a single target task, the transferred policy with our action translator (without conditioning on the task context) yields episode rewards significantly better than the source policy on the target task. Fig. 4 visualizes moving paths of robot arms. The transferred

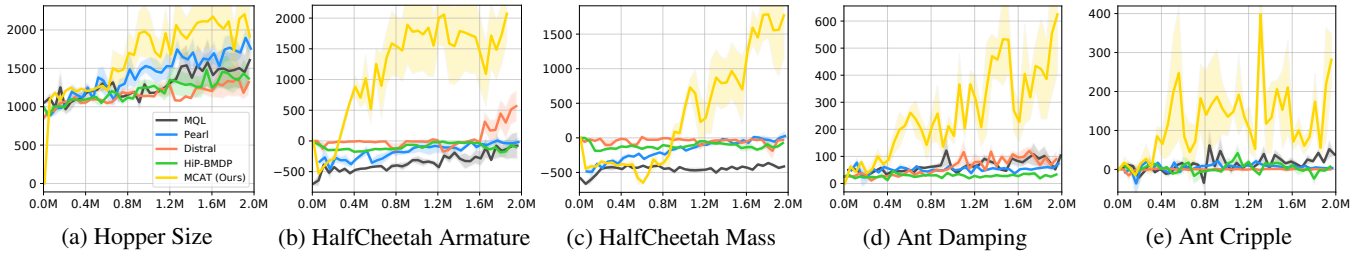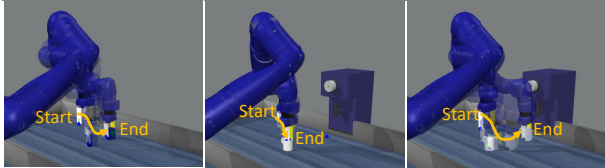| (a) Hopper Size | (b) HalfCheetah Armature | (c) HalfCheetah Mass | (d) Ant Damping | (e) Ant Cripple |

Figure 3: Learning curves of episode rewards on test tasks, averaged over 3 runs. The x-axis is total number of timesteps and the y-axis is average episode reward. Shadow areas indicate standard error.



(a) Source policy on source task  (b) Source policy on target task  (c) Transferred policy on target task

Figure 4: Robot arm moving paths on source (pushing a *cylinder*) or target task (moving a *mug* to a coffee machine).

policy on target task resembles the source policy on source task, while the source policy has trouble grasping the coffee mug on target task. Videos of agents' behavior are in supplementary materials. Tab. 1 reports experimental results of baseline (Zhang et al. 2020c) transferring the source policy based on action correspondence. It proposes to learn an action translator with three loss terms: adversarial loss, domain cycle-consistency loss, and dynamic cycle-consistency loss. Our loss $\mathcal{L}_{trans}$ (Equation 3) draws upon an idea analogous to dynamic cycle-consistency though we have a more expressive forward model $F$ with context variables. When $F$ is strong and reasonably generalizable, domain cycle-consistency loss training the inverse action translator and adversarial loss constraining the distribution of translated action may not be necessary. Ours with a simpler objective function is competitive with Zhang et al. (2020c).



| (a) HalfCheetah | (b) Ant |

Figure 5: Improvement transferred policy over source policy.

We extend the action translator to multiple tasks by conditioning $H$ on context variables of source and target tasks.

We measure the improvement of our transferred policy over the source policy on the target tasks. On HalfCheetah tasks $\mathcal{T}^{(1)} \cdots \mathcal{T}^{(5)}$, the armature becomes larger. As the physics parameter in the target task deviates more from source task, the advantage of transferred policy tends to be more significant (Fig. 5a), because the performance of transferred policy does not drop as much as source policy. We remark that the unified action translator is for any pair of source and target tasks. So action translation for the diagonal elements might be less than $0\%$. For each task on Ant, we set one of its four legs crippled, so any action applied to the crippled leg joints is set as 0. Ideal equivalent action does not always exist across tasks with different crippled legs in this setting. Therefore, it is impossible to minimize $d$ in Proposition 1 as 0. Nevertheless, the inequality proved in Proposition 1 still holds and policy transfer empirically shows positive improvement on most source-target pairs (Fig. 5b).

## 5.2 Comparison with Context-based Meta-RL

We evaluate MCAT combining policy transfer with context-based TD3 in meta-RL problems. The action translator is trained dynamically with data maintained in replay buffer and the source policy keeps being updated. On MuJoCo, we modify environment physics parameters (e.g. size, mass, damping) that affect the transition dynamics to design tasks. We predefine a fixed set of physics parameters for training tasks and unseen test tasks. In order to test algorithms' ability in tackling difficult tasks, environment rewards are delayed to create sparse-reward RL problems (Oh et al. 2018; Tang 2020). In particular, we accumulate dense rewards over $n$ consecutive steps, and the agent receives the delayed feedback every $n$ step or when the episode terminates. To fully exploit the good data collected from our transferred policy, we empirically incorporate self-imitation learning (SIL) (Oh et al. 2018), which imitates the agent's own successful past experiences to further improve the policy.

We compare with several context-based meta-RL methods: MQL (Fakoor et al. 2019), PEARL (Rakelly et al. 2019), Distral (Teh et al. 2017), and HiP-BMDP (Zhang et al. 2020b). Although the baselines perform well on MuJoCo environments with dense rewards, the delayed environment rewards degrade policy learning (Tab. 2, Fig. 3) because the rare transitions with positive rewards are not fully exploited. In contrast, MCAT shows a substantial advantage in performance and sample complexity on both the training tasks and the test tasks. Notably, the perfor-

mance gap is more significant in more complex environments (e.g. HalfCheetah and Ant with higher-dimensional state and sparser rewards). We additionally analyze effect of SIL in Appendix D.4. SIL brings improvements to baselines but MCAT still shows obvious advantages.

| Setting | Hopper Size | Half Cheetah Armature | Half Cheetah Mass | Ant Damp | Ant Cripple |
|---|---|---|---|---|---|
| MQL | 1607.5 | -77.9 | -413.9 | 103.1 | 38.2 |
| PEARL | 1755.8 | -18.8 | 25.9 | 73.2 | 3.5 |
| Distral | 1319.8 | 566.9 | -29.5 | 90.5 | -0.1 |
| HiP-BMDP | 1368.3 | -102.4 | -74.8 | 33.1 | 7.3 |
| MCAT(Ours) | **1914.8** | **2071.5** | **1771.1** | **624.6** | **281.6** |

Table 2: Test rewards at 2M timesteps, averaged over 3 runs.

## 5.3 Ablative Study

**Effect of Policy Transfer** Our MCAT is implemented by combining context-based TD3, self-imitation learning, and policy transfer (PT). We investigate the effect of policy transfer. In Tab. 3. MCAT significantly outperforms MCAT w/o PT, because PT facilitates more balanced performance across training tasks and hence better generalization to test tasks. This empirically confirms that policy transfer is beneficial in meta-RL on sparse-reward tasks.

| Setting | Hopper Size | Half Cheetah Armature | Half Cheetah Mass | Ant Damp | Ant Cripple |
|---|---|---|---|---|---|
| MCAT w/o PT | 1497.5 | 579.1 | -364.3 | 187.7 | 92.4 |
| MCAT | 1982.1 | 1776.8 | 67.1 | 211.8 | 155.7 |
| Improve(%) | 32.3 | 206.8 | 118.4 | 12.8 | 68.5 |

Table 3: Test rewards at 1M timesteps. We report improvements brought by policy transfer (PT).

**Sparser Rewards** We analyze MCAT when rewards are delayed for different numbers of steps (Tab. 4). When rewards are relatively dense (i.e. delay step is 200), during training, the learned policy can reach a high score on each task without the issue of imbalanced performance among multiple tasks. MCAT w/o PT and MCAT perform comparably well within the standard error. However, as the rewards become sparser, it requires longer sequences of correct actions to obtain potentially high rewards. Policy learning struggles on some tasks and policy transfer plays an important role to exploit the precious good experiences on source tasks. Policy transfer brings more improvement on sparser-reward tasks.

In Appendix, we further provide ablative study about More Diverse Tasks (D.3), Effect of SIL (D.4) and Effect of Contrastive Loss (D.5). Appendix D.6 shows that trivially combining the complex action translator (Zhang et al. 2020c) with context-based meta-RL underperforms MCAT.

| Setting | Armature | | | Mass | | |
|---|---|---|---|---|---|---|
| Delay steps | 200 | 350 | 500 | 200 | 350 | 500 |
| MCAT w/o PT | 2583.2 | 1771.7 | 579.1 | 709.6 | 156.6 | -364.2 |
| MCAT | 2251.8 | 2004.5 | 1776.8 | 666.7 | 247.8 | 67.1 |
| Improve(%) | -12.8 | 13.1 | 206.9 | -6.1 | 58.2 | 118.4 |

Table 4: Test rewards at 1M timestpes averaged over 3 runs, on HalfCheetah with *armature* / *mass* changing across tasks.

## 6 Discussion

The scope of MCAT is for tasks with varying dynamics, same as many prior works (Yu et al. 2017; Nagabandi et al. 2018; Nagabandi, Finn, and Levine 2018; Zhou, Pinto, and Gupta 2019). our theory and method of policy transfer can be extended to more general cases (1) tasks with varying reward functions (2) tasks with varying state & action spaces.

Following the idea in Sec. 4, on two general MDPs, we are interested in equivalent state-action pairs achieving the same reward and transiting to equivalent next states. Similar to Proposition 1, we can prove that, on two general MDPs, for two correspondent states $s^{(i)}$ and $s^{(j)}$, the value difference $|V^{\pi^{(i)}}(s^{(i)}, \mathcal{T}^{(i)}) - V^{\pi^{(j)}}(s^{(j)}, \mathcal{T}^{(j)})|$ is upper bounded by $\frac{d}{1-\gamma}$, where $d$ depends on $D_{TV}$ between the next state distribution on source task and the probability distribution of correspondent next state on target task. As an extension, we learn a state translator jointly with our action translator to capture state and action correspondence. Compared with Zhang et al. (2020c) learning both state and action translator, we simplify the objective function training action translator and afford the theoretical foundation. For (1) tasks with varying reward functions, we conduct experiments on Meta-World moving the robot arm to a goal location. The reward at each step is inversely proportional to its distance from the goal location. We fix a goal location on source task and set target tasks with distinct goal locations. Furthermore, we evaluate our approach on 2-leg and 3-leg HalfCheetah. We can test our idea on (2) tasks with varying state and action spaces of different dimensions because the agents have different numbers of joints on the source and target task. Experiments demonstrate that ours with a simpler objective function than the baseline (Zhang et al. 2020c) can transfer the source policy to perform well on the target task. Details of theorems, proofs, and experiments are in Appendix E.

## 7 Conclusion

Meta-RL with long-horizon, sparse-reward tasks is challenging because an agent can rarely obtain positive rewards, and handling multiple tasks simultaneously requires massive samples from distinctive tasks. We propose a simple yet effective objective function to learn an action translator for multiple tasks and provide the theoretical ground. We develop a novel algorithm MCAT using the action translator for policy transfer to improve the performance of off-policy, context-based meta-RL algorithms. We empirically show its efficacy in various environments and verify that our policy transfer can offer substantial gains in sample complexity.

## Acknowledgements

## References

Agarwal, R.; Machado, M. C.; Castro, P. S.; and Bellemare, M. G. 2021. Contrastive Behavioral Similarity Embeddings for Generalization in Reinforcement Learning. *arXiv preprint arXiv:2101.05265*.

Ammar, H. B.; and Taylor, M. E. 2011. Reinforcement learning transfer via common subspaces. In *International Workshop on Adaptive and Learning Agents*, 21–36. Springer.

Duan, Y.; Schulman, J.; Chen, X.; Bartlett, P. L.; Sutskever, I.; and Abbeel, P. 2016. Rl 2: Fast reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*.

Fakoor, R.; Chaudhari, P.; Soatto, S.; and Smola, A. J. 2019. Meta-q-learning. *arXiv preprint arXiv:1910.00125*.

Ferns, N.; Panangaden, P.; and Precup, D. 2004. Metrics for Finite Markov Decision Processes. In *UAI*, volume 4, 162–169.

Finn, C.; Abbeel, P.; and Levine, S. 2017. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning*, 1126–1135. PMLR.

Fujimoto, S.; Hoof, H.; and Meger, D. 2018. Addressing function approximation error in actor-critic methods. In *International Conference on Machine Learning*, 1587–1596. PMLR.

Gupta, A.; Devin, C.; Liu, Y.; Abbeel, P.; and Levine, S. 2017. Learning invariant feature spaces to transfer skills with reinforcement learning. *arXiv preprint arXiv:1703.02949*.

Hadsell, R.; Chopra, S.; and LeCun, Y. 2006. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 2, 1735–1742. IEEE.

Konidaris, G.; and Barto, A. 2006. Autonomous shaping: Knowledge transfer in reinforcement learning. In *Proceedings of the 23rd international conference on Machine learning*, 489–496.

Lee, K.; Seo, Y.; Lee, S.; Lee, H.; and Shin, J. 2020. Context-aware dynamics model for generalization in model-based reinforcement learning. In *International Conference on Machine Learning*, 5757–5766. PMLR.

Li, A. C.; Pinto, L.; and Abbeel, P. 2020. Generalized hindsight for reinforcement learning. *arXiv preprint arXiv:2002.11708*.

Liu, E. Z.; Raghunathan, A.; Liang, P.; and Finn, C. 2020. Explore then Execute: Adapting without Rewards via Factorized Meta-Reinforcement Learning. *arXiv preprint arXiv:2008.02790*.

Mishra, N.; Rohaninejad, M.; Chen, X.; and Abbeel, P. 2017. A simple neural attentive meta-learner. *arXiv preprint arXiv:1707.03141*.

Mnih, V.; Kavukcuoglu, K.; Silver, D.; Rusu, A. A.; Veness, J.; Bellemare, M. G.; Graves, A.; Riedmiller, M.; Fidjeland, A. K.; Ostrovski, G.; Petersen, S.; Beattie, C.; Sadik, A.; Antonoglou, I.; King, H.; Kumaran, D.; Wierstra, D.; Legg, S.; and Hassabis, D. 2015. Human-level control through deep reinforcement learning. *Nature*.

Nagabandi, A.; Clavera, I.; Liu, S.; Fearing, R. S.; Abbeel, P.; Levine, S.; and Finn, C. 2018. Learning to adapt in dynamic, real-world environments through meta-reinforcement learning. *arXiv preprint arXiv:1803.11347*.

Nagabandi, A.; Finn, C.; and Levine, S. 2018. Deep online learning via meta-learning: Continual adaptation for model-based rl. *arXiv preprint arXiv:1812.07671*.

Oh, J.; Guo, Y.; Singh, S.; and Lee, H. 2018. Self-imitation learning. In *International Conference on Machine Learning*, 3878–3887. PMLR.

Parisotto, E.; Ba, J. L.; and Salakhutdinov, R. 2015. Actor-mimic: Deep multitask and transfer reinforcement learning. *arXiv preprint arXiv:1511.06342*.

Rakelly, K.; Zhou, A.; Finn, C.; Levine, S.; and Quillen, D. 2019. Efficient off-policy meta-reinforcement learning via probabilistic context variables. In *International conference on machine learning*, 5331–5340. PMLR.

Ren, H.; Zhu, Y.; Leskovec, J.; Anandkumar, A.; and Garg, A. 2020. OCEAN: Online Task Inference for Compositional Tasks with Context Adaptation. In *Conference on Uncertainty in Artificial Intelligence*, 1378–1387. PMLR.

Rusu, A. A.; Colmenarejo, S. G.; Gulcehre, C.; Desjardins, G.; Kirkpatrick, J.; Pascanu, R.; Mnih, V.; Kavukcuoglu, K.; and Hadsell, R. 2015. Policy distillation. *arXiv preprint arXiv:1511.06295*.

Schulman, J.; Wolski, F.; Dhariwal, P.; Radford, A.; and Klimov, O. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.

Seo, Y.; Lee, K.; Clavera, I.; Kurutach, T.; Shin, J.; and Abbeel, P. 2020. Trajectory-wise Multiple Choice Learning for Dynamics Generalization in Reinforcement Learning. *arXiv preprint arXiv:2010.13303*.

Silver, D.; Huang, A.; Maddison, C. J.; Guez, A.; Sifre, L.; Van Den Driessche, G.; Schrittwieser, J.; Antonoglou, I.; Panneershelvam, V.; Lanctot, M.; et al. 2016. Mastering the game of Go with deep neural networks and tree search. *nature*, 529(7587): 484.

Stadie, B. C.; Yang, G.; Houthooft, R.; Chen, X.; Duan, Y.; Wu, Y.; Abbeel, P.; and Sutskever, I. 2018. Some considerations on learning to explore via meta-reinforcement learning. *arXiv preprint arXiv:1803.01118*.

Sung, F.; Zhang, L.; Xiang, T.; Hospedales, T.; and Yang, Y. 2017. Learning to learn: Meta-critic networks for sample efficient learning. *arXiv preprint arXiv:1706.09529*.

Tang, Y. 2020. Self-imitation learning via generalized lower bound q-learning. *arXiv preprint arXiv:2006.07442*.

Teh, Y. W.; Bapst, V.; Czarnecki, W. M.; Quan, J.; Kirkpatrick, J.; Hadsell, R.; Heess, N.; and Pascanu, R. 2017. Distral: Robust multitask reinforcement learning. *arXiv preprint arXiv:1707.04175*.

Todorov, E.; Erez, T.; and Tassa, Y. 2012. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 5026–5033. IEEE.

Van der Maaten, L.; and Hinton, G. 2008. Visualizing data using t-SNE. *Journal of machine learning research*, 9(11).

Wang, J. X.; Kurth-Nelson, Z.; Tirumala, D.; Soyer, H.; Leibo, J. Z.; Munos, R.; Blundell, C.; Kumaran, D.; and Botvinick, M. 2016. Learning to reinforcement learn. *arXiv preprint arXiv:1611.05763*.

Xu, Z.; van Hasselt, H.; and Silver, D. 2018. Meta-gradient reinforcement learning. *arXiv preprint arXiv:1805.09801*.

Yin, H.; and Pan, S. 2017. Knowledge transfer for deep reinforcement learning with hierarchical experience replay. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 1.

Yu, T.; Kumar, S.; Gupta, A.; Levine, S.; Hausman, K.; and Finn, C. 2020a. Gradient surgery for multi-task learning. *arXiv preprint arXiv:2001.06782*.

Yu, T.; Quillen, D.; He, Z.; Julian, R.; Hausman, K.; Finn, C.; and Levine, S. 2020b. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on Robot Learning*, 1094–1100. PMLR.

Yu, W.; Tan, J.; Liu, C. K.; and Turk, G. 2017. Preparing for the unknown: Learning a universal policy with online system identification. *arXiv preprint arXiv:1702.02453*.

Zhang, A.; Lyle, C.; Sodhani, S.; Filos, A.; Kwiatkowska, M.; Pineau, J.; Gal, Y.; and Precup, D. 2020a. Invariant causal prediction for block mdps. In *International Conference on Machine Learning*, 11214–11224. PMLR.

Zhang, A.; Sodhani, S.; Khetarpal, K.; and Pineau, J. 2020b. Learning robust state abstractions for hidden-parameter block {mdp} s. In *International Conference on Learning Representations*.

Zhang, Q.; Xiao, T.; Efros, A. A.; Pinto, L.; and Wang, X. 2020c. Learning Cross-Domain Correspondence for Control with Dynamics Cycle-Consistency. *arXiv preprint arXiv:2012.09811*.

Zhou, W.; Pinto, L.; and Gupta, A. 2019. Environment probing interaction policies. *arXiv preprint arXiv:1907.11740*.

Zhu, Z.; Lin, K.; and Zhou, J. 2020. Transfer Learning in Deep Reinforcement Learning: A Survey. *arXiv preprint arXiv:2009.07888*.

*Appendix:*

# A   Bound Value Difference in Policy Transfer

In this section, we provide detailed theoretical ground for our policy transfer approach, as a supplement to Sec. 4. We first define a binary relation for actions to describe the correspondent actions behaving equivalently on two MDPs (Definition 1). Building upon the notion of action equivalence, we derive the upper bound of value difference between policies on two MDPs (Theorem 1). Finally, we reach a proposition for the upper bound of value difference (Proposition 1) to explain that minimizing our objective function results in bounding the value difference between the source and transferred policy.

**Definition 1.** *Given two MDPs $\mathcal{T}^{(i)} = \{\mathcal{S}, \mathcal{A}, p^{(i)}, r^{(i)}, \gamma, \rho_0\}$ and $\mathcal{T}^{(j)} = \{\mathcal{S}, \mathcal{A}, p^{(j)}, r^{(j)}, \gamma, \rho_0\}$ with the same state space and action space, for each state $s \in \mathcal{S}$, we define a binary relation $B_s \in \mathcal{A} \times \mathcal{A}$ called **action equivalence relation**. For any action $a^{(i)} \in \mathcal{A}$, $a^{(j)} \in \mathcal{A}$, if $(a^{(i)}, a^{(j)}) \in B_s$ (i.e. $a^{(i)} B_s a^{(j)}$), the following conditions hold:*

$$r^{(i)}(s, a^{(i)}) = r^{(j)}(s, a^{(j)}) \text{ and } p^{(i)}(\cdot|s, a^{(i)}) = p^{(j)}(\cdot|s, a^{(j)}) \tag{5}$$

Based on Definition 1, at state $s$, action $a^{(i)}$ on $\mathcal{T}^{(i)}$ is equivalent to action $a^{(j)}$ on $\mathcal{T}^{(j)}$ if $a^{(i)} B_s a^{(j)}$. Note that the binary relation $B_s$ is defined for each $s$ separately. The action equivalence relation might change on varied states. On two MDPs with the same dynamic and reward functions, it is trivial to get the equivalent action with identity mapping. However, we are interested in more complex cases where the reward and dynamic functions are not identical on two MDPs.

Ideally, the equivalent action always exists on the target MDP $\mathcal{T}^{(i)}$ for any state-action pair on the source MDP $\mathcal{T}^{(j)}$ and there exists an action translator function $H : \mathcal{S} \times \mathcal{A} \to \mathcal{A}$ to identify the exact equivalent action. Starting from state $s$, the translated action $\tilde{a} = H(s, a)$ on the task $\mathcal{T}^{(i)}$ generates reward and next state distribution the same as action $a$ on the task $\mathcal{T}^{(j)}$ (i.e. $\tilde{a} B_s a$). Then any deterministic policy $\pi^{(j)}$ on the source task $\mathcal{T}^{(j)}$ can be perfectly transferred to the target task $\mathcal{T}^{(i)}$ with $\pi^{(i)}(s) = H(s, \pi^{(j)}(s))$. The value of the policy $\pi^{(j)}$ on the source task $\mathcal{T}^{(j)}$ is equal to the value of transferred policy $\pi^{(i)}$ on the target task $\mathcal{T}^{(i)}$.

Without the assumption of existence of a perfect correspondence for each action, given any two deterministic policies $\pi^{(j)}$ and $\pi^{(i)}$, we prove that the difference in the policy value is upper bounded by a scalar $\frac{d}{1-\gamma}$ depending on L1-distance between reward functions $|r^{(j)}(s, \pi^{(j)}(s)) - r^{(i)}(s, \pi^{(i)}(s))|$ and total-variation distance between next state distributions $D_{TV}(p^{(j)}(\cdot|s, \pi^{(j)}(s)), p^{(i)}(\cdot|s, \pi^{(i)}(s)))$.

**Theorem 1.** *Let $\mathcal{T}^{(i)} = \{\mathcal{S}, \mathcal{A}, p^{(i)}, r^{(i)}, \gamma, \rho_0\}$ and $\mathcal{T}^{(j)} = \{\mathcal{S}, \mathcal{A}, p^{(j)}, r^{(j)}, \gamma, \rho_0\}$ be two MDPs sampled from the distribution of tasks $p(\mathcal{T})$. $\pi^{(i)}$ is a deterministic policy on $\mathcal{T}^{(i)}$ and $\pi^{(j)}$ is a deterministic policy on $\mathcal{T}^{(j)}$. Let $M = \sup_{s \in \mathcal{S}} |V^{\pi^{(j)}}(s, \mathcal{T}^{(j)})|$, $d = \sup_{s \in \mathcal{S}} \left[ |r^{(j)}(s, \pi^{(j)}(s)) - r^{(i)}(s, \pi^{(i)}(s))| + 2\gamma M D_{TV}(p^{(j)}(\cdot|s, \pi^{(j)}(s)), p^{(i)}(\cdot|s, \pi^{(i)}(s))) \right]$. For $\forall s \in \mathcal{S}$, we have*

$$\left| V^{\pi^{(i)}}(s, \mathcal{T}^{(i)}) - V^{\pi^{(j)}}(s, \mathcal{T}^{(j)}) \right| \leq \frac{d}{1-\gamma} \tag{6}$$

*Proof.* Let $a^{(i)} = \pi^{(i)}(s)$ and $a^{(j)} = \pi^{(j)}(s)$. $s'$ denotes the next state following state $s$. $s''$ denotes the next state following $s'$. We rewrite the value difference as:

$$
\begin{aligned}
V^{\pi^{(j)}}(s, \mathcal{T}^{(j)}) - V^{\pi^{(i)}}(s, \mathcal{T}^{(i)}) &= r^{(j)}(s, a^{(j)}) + \gamma \sum_{s'} p^{(j)}(s'|s, a^{(j)}) V^{\pi^{(j)}}(s', \mathcal{T}^{(j)}) \\
&\quad - r^{(i)}(s, a^{(i)}) - \gamma \sum_{s'} p^{(i)}(s'|s, a^{(i)}) V^{\pi^{(i)}}(s', \mathcal{T}^{(i)}) \\
&= (r^{(j)}(s, a^{(j)}) - r^{(i)}(s, a^{(i)})) \\
&\quad + \gamma \left[ \sum_{s'} p^{(j)}(s'|s, a^{(j)}) V^{\pi^{(j)}}(s', \mathcal{T}^{(j)}) - \sum_{s'} p^{(i)}(s'|s, a^{(i)}) V^{\pi^{(i)}}(s', \mathcal{T}^{(i)}) \right] \\
&\qquad \text{\textcolor{blue}{*minus and plus} } \gamma \sum_{s'} p^{(i)}(s'|s, a^{(i)}) V^{\pi^{(j)}}(s', \mathcal{T}^{(j)}) \\
&= (r^{(j)}(s, a^{(j)}) - r^{(i)}(s, a^{(i)})) \\
&\quad + \gamma \sum_{s'} \left[ p^{(j)}(s'|s, a^{(j)}) - p^{(i)}(s'|s, a^{(i)}) \right] V^{\pi^{(j)}}(s', \mathcal{T}^{(j)}) \\
&\quad + \gamma \sum_{s'} p^{(i)}(s'|s, a^{(i)}) \left[ V^{\pi^{(j)}}(s', \mathcal{T}^{(j)}) - V^{\pi^{(i)}}(s', \mathcal{T}^{(i)}) \right]
\end{aligned}
$$

Then we consider the absolute value of the value difference:

$$
\begin{aligned}
\left| V^{\pi^{(j)}}(s, \mathcal{T}^{(j)}) - V^{\pi^{(i)}}(s, \mathcal{T}^{(i)}) \right| &\leq \left| r^{(j)}(s, a^{(j)}) - r^{(i)}(s, a^{(i)}) \right| \\
&+ \gamma \left| \sum_{s'} \left[ p^{(j)}(s'|s, a^{(j)}) - p^{(i)}(s'|s, a^{(i)}) \right] V^{\pi^{(j)}}(s', \mathcal{T}^{(j)}) \right| \\
&+ \gamma \left| \sum_{s'} p^{(i)}(s'|s, a^{(i)}) \left[ V^{\pi^{(j)}}(s', \mathcal{T}^{(j)}) - V^{\pi^{(i)}}(s', \mathcal{T}^{(i)}) \right] \right| \\
&\textcolor{blue}{\text{*property of total variation distance when the set is countable}} \\
&= \left| r^{(j)}(s, a^{(j)}) - r^{(i)}(s, a^{(i)}) \right| \\
&+ 2\gamma M D_{TV}(p^{(j)}(\cdot|s, a^{(j)}), p^{(i)}(\cdot|s, a^{(i)})) \\
&+ \gamma \left| \sum_{s'} p^{(i)}(s'|s, a^{(i)}) \left[ V^{\pi^{(j)}}(s', \mathcal{T}^{(j)}) - V^{\pi^{(i)}}(s', \mathcal{T}^{(i)}) \right] \right| \\
&\leq d + \gamma \left| \sum_{s'} p^{(i)}(s'|s, a^{(i)}) \left[ V^{\pi^{(j)}}(s', \mathcal{T}^{(j)}) - V^{\pi^{(i)}}(s', \mathcal{T}^{(i)}) \right] \right| \\
&\leq d + \gamma \sup_{s'} \left| V^{\pi^{(j)}}(s', \mathcal{T}^{(j)}) - V^{\pi^{(i)}}(s', \mathcal{T}^{(i)}) \right| \\
&\textcolor{blue}{\text{*by induction}} \\
&\leq d + \gamma \left[ d + \gamma \sup_{s''} \left| V^{\pi^{(j)}}(s'', \mathcal{T}^{(j)}) - V^{\pi^{(i)}}(s'', \mathcal{T}^{(i)}) \right| \right] \\
&\leq d + \gamma d + \gamma^2 \sup_{s''} \left| V^{\pi^{(j)}}(s'', \mathcal{T}^{(j)}) - V^{\pi^{(i)}}(s'', \mathcal{T}^{(i)}) \right| \\
&\leq \cdots \\
&\leq d + \gamma d + \gamma^2 d + \gamma^3 d + \cdots = \frac{d}{1 - \gamma}
\end{aligned}
$$

$\square$

For a special case where reward function $r(s, a, s')$ only depends on the current state $s$ and next state $s'$, the upper bound of policy value difference is only related to the distance in next state distributions.

**Proposition 1.** *Let $\mathcal{T}^{(i)} = \{\mathcal{S}, \mathcal{A}, p^{(i)}, r^{(i)}, \gamma, \rho_0\}$ and $\mathcal{T}^{(j)} = \{\mathcal{S}, \mathcal{A}, p^{(j)}, r^{(j)}, \gamma, \rho_0\}$ be two MDPs sampled from the distribution of tasks $p(\mathcal{T})$. $\pi^{(i)}, \pi^{(j)}$ is the deterministic policy on $\mathcal{T}^{(i)}, \mathcal{T}^{(j)}$. Assume the reward function only depends on the state and next state $r^{(i)}(s, a^{(i)}, s') = r^{(j)}(s, a^{(j)}, s') = r(s, s')$. Let $d = \sup_{s \in \mathcal{S}} 2M D_{TV}(p^{(j)}(\cdot|s, \pi^{(j)}(s)), p^{(i)}(\cdot|s, \pi^{(i)}(s)))$ and $M = \sup_{s \in \mathcal{S}, s' \in \mathcal{S}} |r(s, s') + \gamma V^{\pi^{(j)}}(s, \mathcal{T}^{(j)})|$. $\forall s \in \mathcal{S}$, we have*

$$
\left| V^{\pi^{(i)}}(s, \mathcal{T}^{(i)}) - V^{\pi^{(j)}}(s, \mathcal{T}^{(j)}) \right| \leq \frac{d}{1 - \gamma} \tag{4}
$$

*Proof.* Let $a^{(i)} = \pi^{(i)}(s)$ and $a^{(j)} = \pi^{(j)}(s)$. $s'$ denotes the next state following state $s$. $s''$ denotes the next state following $s'$.
In the special case of $r^{(i)}(s, a^{(i)}, s') = r(s, s')$, the value of policy can be written as:

$$
\begin{aligned}
V^{\pi^{(i)}}(s, \mathcal{T}^{(i)}) &= r^{(i)}(s, a^{(i)}) + \gamma \sum_{s'} p^{(i)}(s'|s, a^{(i)}) V^{\pi^{(i)}}(s', \mathcal{T}^{(i)}) \\
&= \sum_{s'} p^{(i)}(s'|s, a^{(i)}) r(s, s') + \gamma \sum_{s'} p^{(i)}(s'|s, a^{(i)}) V^{\pi^{(i)}}(s', \mathcal{T}^{(i)}) \\
&= \sum_{s'} p^{(i)}(s'|s, a^{(i)}) \left[ r(s, s') + \gamma V^{\pi^{(i)}}(s', \mathcal{T}^{(i)}) \right]
\end{aligned}
$$

We can derive the value difference:

$$V^{\pi^{(j)}}(s, \mathcal{T}^{(j)}) - V^{\pi^{(i)}}(s, \mathcal{T}^{(i)}) \;=\; \sum_{s'} p^{(j)}(s'|s, a^{(j)}) \left[ r(s, s') + \gamma V^{\pi^{(j)}}(s', \mathcal{T}^{(j)}) \right] - \sum_{s'} p^{(i)}(s'|s, a^{(i)}) \left[ r(s, s') + \gamma V^{\pi^{(i)}}(s', \mathcal{T}^{(i)}) \right]$$

<span style="color:blue">*minus and plus $\sum_{s'} p^{(i)}(s'|s, a^{(i)}) \left[ r(s, s') + \gamma V^{\pi^{(j)}}(s', \mathcal{T}^{(j)}) \right]$</span>

<span style="color:blue">**combine the first two terms, combine the last two terms</span>

$$\;=\; \sum_{s'} \left[ p^{(j)}(s'|s, a^{(j)}) - p^{(i)}(s'|s, a^{(i)}) \right] \left[ r(s, s') + \gamma V^{\pi^{(j)}}(s', \mathcal{T}^{(j)}) \right]$$
$$+\; \gamma \sum_{s'} p^{(i)}(s'|s, a^{(i)}) \left[ V^{\pi^{(j)}}(s', \mathcal{T}^{(j)}) - V^{\pi^{(i)}}(s', \mathcal{T}^{(i)}) \right]$$

Then we take absolute value of the value difference:

$$\left| V^{\pi^{(j)}}(s, \mathcal{T}^{(j)}) - V^{\pi^{(i)}}(s, \mathcal{T}^{(i)}) \right| \;\leq\; 2M D_{TV}\big(p^{(j)}(\cdot|s, a^{(j)}), p^{(i)}(\cdot|s, a^{(i)})\big)$$
$$+\; \gamma \left| \sum_{s'} p^{(i)}(s'|s, a^{(i)}) \left[ V^{\pi^{(j)}}(s', \mathcal{T}^{(j)}) - V^{\pi^{(i)}}(s', \mathcal{T}^{(i)}) \right] \right|$$
$$\leq\; d + \gamma \left| \sum_{s'} p^{(i)}(s'|s, a^{(i)}) \left[ V^{\pi^{(j)}}(s', \mathcal{T}^{(j)}) - V^{\pi^{(i)}}(s', \mathcal{T}^{(i)}) \right] \right|$$
$$\leq\; d + \gamma \sup_{s'} \left| V^{\pi^{(j)}}(s', \mathcal{T}^{(j)}) - V^{\pi^{(i)}}(s', \mathcal{T}^{(i)}) \right|$$
$$\leq\; d + \gamma \left[ d + \gamma \sup_{s''} \left| V^{\pi^{(j)}}(s'', \mathcal{T}^{(j)}) - V^{\pi^{(i)}}(s'', \mathcal{T}^{(i)}) \right| \right]$$
$$\leq\; d + \gamma d + \gamma^2 \sup_{s''} \left| V^{\pi^{(j)}}(s'', \mathcal{T}^{(j)}) - V^{\pi^{(i)}}(s'', \mathcal{T}^{(i)}) \right|$$
$$\leq\; \cdots$$
$$\leq\; d + \gamma d + \gamma^2 d + \gamma^3 d + \cdots = \frac{d}{1-\gamma}$$

$\square$

# B   Algorithm of MCAT

---

**Algorithm 1: MCAT combining context-based meta-RL algorithm with policy transfer**

---

1: Initialize critic networks $Q_{\theta_1}$, $Q_{\theta_2}$ and actor network $\pi_\phi$ with random parameters $\theta_1$, $\theta_2$, $\phi$
2: Initialize target networks $\theta_1' \leftarrow \theta_1$, $\theta_2' \leftarrow \theta_2$, $\phi' \leftarrow \phi$
3: Initialize replay buffer $\mathcal{B} = \mathcal{B}^{(1)} \cup \mathcal{B}^{(2)} \cup \cdots \cup \mathcal{B}^{(|\mathcal{T}|)}$ and $\mathcal{B}^{(i)} \leftarrow \emptyset$ for each $i$.
4: Initialize SIL replay buffer $\mathcal{D} \leftarrow \emptyset$
5: Initialize context encoder $C_{\psi_C}$, forward model $F_{\psi_F}$, action translator $H_{\psi_H}$
6: Initialize set of trajectory rewards for shared policy on each task in recent timesteps as $R^{(i)} = \emptyset$, set of trajectory rewards for transferred policy from $\mathcal{T}^{(j)}$ to $\mathcal{T}^{(i)}$ in recent timesteps as $R^{(j)\rightarrow(i)} = \emptyset$. $\bar{R}$ denotes average episode rewards in the set.
7: **for** each iteration **do**
8:    // Collect training samples
9:    **for** each task $\mathcal{T}^{(i)}$ **do**
10:       **if** $R^{(i)} = \emptyset$ **then**
11:          use the shared policy in this episode
12:       **else if** there exist $j \in 1, 2, \cdots, |\mathcal{T}|$ such that $R^{(j)\rightarrow(i)} = \emptyset$ and $\bar{R}^{(j)} > \bar{R}^{(i)}$ **then**
13:          use transferred policy from source task $\mathcal{T}^{(j)}$ to target task $\mathcal{T}^{(i)}$ in this episode
14:       **else if** there exist $j \in 1, 2, \cdots, |\mathcal{T}|$, such that $j = \arg\max_{j'} \bar{R}^{(j')\rightarrow(i)}$ and $\bar{R}^{(j)\rightarrow(i)} > \bar{R}^{(i)}$ **then**
15:          use transferred policy from source task $\mathcal{T}^{(j)}$ to target task $\mathcal{T}^{(i)}$ in this episode
16:       **else**
17:          use the shared policy in this episode
18:       **end if**
19:       **for** $t = 1$ to TaskHorizon **do**
20:          Get context latent variable $z_t = C_{\psi_C}(\tau_{t,K})$
21:          Select the action $a$ based on the transferred policy or shared policy, take the action with noise $a_t = a + \epsilon$ where $\epsilon \sim \mathcal{N}(0, \sigma)$, observe reward $r_t$ and new state $s_{t+1}$. Update $\mathcal{B}^{(i)} \leftarrow \mathcal{B}^{(i)} \cup \{s_t, a_t, r_t, s_{t+1}, \tau_{t,K}\}$
22:       **end for**
23:       Compute returns $R_t = \sum_{k=t}^{\infty} \gamma^{k-t} r_k$ and update $\mathcal{D} \leftarrow \mathcal{D} \cup \{s_t, a_t, r_t, s_{t+1}, \tau_{t,K}, R_t\}$ for every step $t$ in this episode.

24:       Update the average reward of shared policy on task $\mathcal{T}^{(i)}$ (i.e. $R^{(i)}$) if we took shared policy in this episode, or update the average reward of the transferred policy from $\mathcal{T}^{(j)}$ to $\mathcal{T}^{(i)}$ (i.e. $R^{(j)\rightarrow(i)}$) if we took the transferred policy.
25:    **end for**
26:    // Update the context encoder $C_{\psi_C}$ and forward model $F_{\psi_F}$ with $\mathcal{L}_{forw}$ and $\mathcal{L}_{cont}$
27:    // Update the action translator $H_{\psi_H}$ with $\mathcal{L}_{trans}$
28:    // Update the critic network $Q_{\theta_1}$, $Q_{\theta_2}$ and actor network $\pi_\phi$ with TD3 and SIL objective function
29:    **for** step in training steps **do**
30:       Update $\theta_1$, $\theta_2$ for the critic networks to minimize $\mathcal{L}_{td3} + \mathcal{L}_{sil}$ (see Algorithm 2)
31:       Update $\phi$ for the actor network with deterministic policy gradient
32:       Update the $\theta_1'$, $\theta_2'$, $\phi'$ for target networks with soft assignment
33:    **end for**
34:    // Update the trajectory reward for shared policy and transferred policy if necessary
35:    **for** each task $\mathcal{T}^{(i)}$ **do**
36:       pop out trajectory rewards in $R^{(i)}$ which were stored before the last G timesteps
37:       pop out trajectory rewards in $R^{(j)\rightarrow(i)}(\forall j)$ which were stored before the last G timesteps
38:    **end for**
39: **end for**

---

**Algorithm 2:** Compute critic loss based on TD3 algorithm and SIL algorithm

---
1: Sample batch data of transitions $(s_t, a_t, r_t, s_{t+1}, \tau_{t,K}) \in \mathcal{B}$
2: Get context variable $z_t = C_{\psi_C}(\tau_{t,K})$. Get next action $a_{t+1} \sim \pi_{\phi'}(z_t, s_{t+1}) + \epsilon, \epsilon \sim \text{clip}(\mathcal{N}(0, \tilde{\sigma}), -c, c)$
3: Get target value for critic network $y = r_t + \gamma \min_{l=1,2} Q_{\theta_l'}(z_t, s_{t+1}, a_{t+1})$.
4: Compute TD error $\mathcal{L}_{td3} = \min_{l=1,2}(y - Q_{\theta_l}(z_t, s_{t+1}, a_{t+1}))^2$
5: Sample batch data of transitions $(s_t, a_t, \tau_{t,K}, R_t) \in \mathcal{D}$
6: Get context variable $z_t = C_{\psi_C}(\tau_{t,K})$. Compute SIL loss $\mathcal{L}_{sil} = \sum_{l=1,2} \max(R_t - Q_{\theta_l}(z_t, s_t, a_t), 0)^2$

---

# C    Experiment Details

In this section, we explain more details for Section 5 and show additional experimental results.

## C.1    Environment

**MuJoCo**    We use Hopper, HalfCheetah and Ant environments from OpenAI Gym based on the MuJoCo physics engine (Todorov, Erez, and Tassa 2012). The goal is to move forward while keeping the control cost minimal.

- **Hopper** Hopper agent consists of 5 rigid links with 3 joints. Observation $s_t$ is an 11-dimension vector consisting of root joint's position (except for x-coordinate) and velocity angular position and velocity of all 3 joints. Action $a_t$ lies in the space $[-1.0, 1.0]^3$, which corresponds to the torques applied to 3 joints. Reward $r_t = v_{\text{torso},t} - 0.001\|a_t\|^2 + 1.0$ means the forward velocity of the torso $v_{\text{torso},t}$ minus the control cost for action $0.001\|a_t\|^2$ and plus the survival bonus 1.0 at each step. We modify the size of each rigid part to enlarge/contract the body of the agent, so we can create tasks with various dynamics.

- **HalfCheetah** Half-cheetah agent consists of of 7 rigid links (1 for torso, 3 for forelimb, and 3 for hindlimb), connected by 6 joints. State $s_t$ is a 17-dimension vector consisting of root joint's position (except for x-coordinate) and velocity, angular position and velocity of all 6 joints. Action $a_t$ is sampled from the space $[-1.0, 1.0]^6$, representing the torques applied to each of the 6 joints. Reward $r_t = v_{\text{torso},t} - 0.1\|a_t\|^2$ is the forward velocity of the torso minus the control cost for action. In order to design multiple tasks with varying dynamics on HalfCheetah, we modify the armature value (similarly to (Zhang et al. 2020c)) or scale the mass of each rigid link by a fixed scale factor (similarly to (Lee et al. 2020)).

- **Ant** Ant agent consists of 13 rigid links connected by 8 joints. Observation $s_t$ is a 27-dimension vector including information about the root joint's position and velocity, angular position and velocity of all 8 joints, and frame orientations. Action $a_t \in [-1.0, 1.0]^8$ is the torques applied to each of 8 joints. Reward is $r_t = v_{\text{torso},t} + 0.05$, meaning the velocity of moving forward plus the survival bonus 0.05 for each step. To change the environment dynamics, we modify the damping of every leg. Specifically, given a scale factor $d$, we modify two legs to have damping multiplied by $d$, and the other two legs to have damping multiplied by $1/d$ (similarly to (Lee et al. 2020)). Alternatively, we can cripple one of the agent's four legs to change the dynamics function. The torques applied to two joints on the crippled leg (i.e. two correspondent elements in actions) are set as 0. (similarly to (Seo et al. 2020)).

**MetaWorld**    Additionally, we consider the tasks of pushing Cylinder, Coffee Mug and Cube. They are named as push-v2, coffee-push-v2, and sweep-into-goal-v2 on MetaWorld benchmark (Yu et al. 2020b) respectively. The goal is to move the objects from a random initial location to a random goal location. The observation is of dimension 14, consisting of the location of the robot hand, the distance between two gripper fingers, the location and position of the object, and the target location. The action $a \in [-1.0, 1.0]^4$ controls the movement of the robot hand and opening/closing of the gripper. The reward is 1.0 when the object is close to the target location (i.e. distance less than 0.05). Otherwise, the environment reward is 0.0. The length of an episode is 500 steps. The tasks of manipulating different objects have different dynamics. We change the physics parameters armature and damping across tasks to make the policy transfer more challenging.

## C.2    Implementation Details for Policy Transfer with Fixed Dataset & Source Policy

In Section 5.1, we study the performance of policy transfer with our action translator with a fixed dataset and source policy. In this experiment, we demonstrate our proposed policy transfer approach trained with fixed datasets and source policy outperforms the baselines. We provide the experimental details as follows.

**Source Policy and Dataset**

- **MuJoCo** On HalfCheetah, the armature value on the source and target task is 0.1 and 0.5 respectively. On Ant, the leg 0 is crippled on the source task while the leg 3 is crippled on the target task. We train well-performing policies on the source tasks as source policies, and we also train mediocre policies on both source tasks and target tasks to obtain training data.

  We apply the TD3 algorithm(Fujimoto, Hoof, and Meger 2018) and dense rewards to learn policies. The hyperparameters for the TD3 algorithm are listed in Table 5. Specifically, during the start 25K timesteps, the TD3 agent collects data by randomly sampling from the action space. After the first 25K timesteps, the agent learns an deterministic policy based on the data collected in the replay buffer. During training, the agent collects data with actions following the learned policy with Gaussian noise, and updates the replay buffer as well. On HalfCheetah environment, we use the learned policy at 300K timesteps as good policy, and use the learned policy at 80K timesteps as mediocre policy. On Ant environment, the learned policy at 400K timesteps and 20K timesteps are used as good policy and mediocre policy respectively.

  With the mediocre policies, we collect 100K transition samples on the source and target tasks respectively. During data collection, at each step, we record the following information: (a) current state; (b) current action drawn from the mediocre policies; (c) next state; (d) historical observations in the past 10 steps; (e) historical actions in the past 10 steps. The historical transition information are employed to learn the context model for forward dynamics prediction.

| Parameter name | Value |
|---|---|
| Start Timesteps | 2.5e4 |
| Gaussian exploration noise $\sigma$ | 0.1 |
| Batch Size | 256 |
| Discount $\gamma$ | 0.99 |
| Target network update rate | 5e-3 |
| Policy noise $\tilde{\sigma}$ | 0.2 |
| Noise clip $c$ | 0.5 |
| Policy update frequency | 2 |
| Replay buffer size | 1e6 |
| Actor learning rate | 3e-4 |
| Critic learning rate | 3e-4 |
| Optimizer | Adam |
| Actor layers | 3 |
| Hidden dimension | 256 |

Table 5: The hyperparameters for TD3 algorithm.

- **MetaWorld** On source tasks, we keep the default physics parameters. However, on the target task,the value of armature and damping for the gripper joints is 0.1 multiplying the default. We get the manually designed good policies from official public code[1]. The performance of the good source policy is shown in Tab. 6. By adding Gaussian noise following $\mathcal{N}(0, 1.0)$ to action drawn from the good policies, we collect 100K transition samples on the source and target tasks respectively.

With the fixed datasets on both source and targe tasks, we can train action translator to transfer the fixed source policy. First, we learn the forward dynamics model. Then we learn the action translator based on the well-trained forward dynamics model. For fair comparison, we train the baseline (Zhang et al. 2020c) and our action translator with the same dataset and source policy. The hyperparameters and network structures applied in the baseline and our approach are introduced as follows

**Transferred Policy (Zhang et al. 2020c)** This baseline is implemented using the code provided by Zhang et al. (2020c) [2]. The forward dynamics model first encodes the state and action as 128-dimensional vectors respectively via a linear layer with ReLU activation. The state embedding and action embedding is then concatenated to predict the next state with an MLP with 2 hidden layers of 256 units and ReLU activation. We train the forward dynamics model with batch size 32 and decaying learning rate from 0.001, 0.0003 to 0.0001. In order to optimize the forward dynamics model, the objective function is L1-loss between the predicted next state and the actual next state. With these hyper-parameters settings, we train the forward modelFand the context modelCfor30 epochs, each epoch with 10K steps.

The action translator first encodes the state and action as 128-dimensional vectors respectively via a linear layer with ReLU activation. The state embedding and action embedding are then concatenated to generate the translated action via an MLP with 2 hidden layers of 256 units and ReLU activation. As for the objective function with three terms: adversarial loss, domain cycle-consistency loss, and dynamics cycle-consistency loss, we tune three weights. We train the action translator for 30 epochs. After each epoch, the performance of transferred policy with the action translator is evaluated on the target task. We average episode rewards in 100 episodes as the epoch performance. Finally, we report the best epoch performance over the 30 epochs.

| Setting | Source policy on source task | Source policy on target task | Transferred policy (Zhang et al. 2020c) on target task | Transferred policy (Ours) on target task |
|---|---|---|---|---|
| HalfCheetah | 5121.4 | 2355.0 | **3017.1**$_{(\pm44.2)}$ | 2937.2$_{(\pm9.5)}$ |
| Ant | 476.8 | 55.8 | 97.2$_{(\pm2.5)}$ | **208.1**$_{(\pm8.2)}$ |
| Cylinder-Mug | 317.3 | 0.0 | 308.1$_{(\pm75.3)}$ | **395.6**$_{(\pm19.4)}$ |
| Cylinder-Cube | 439.7 | 0.0 | 262.4$_{(\pm48.1)}$ | **446.1**$_{(\pm1.1)}$ |

Table 6: Performance of source and transferred policy on target task. This is expanding Tab. 1 in the main text.

---

[1]https://github.com/rlworkgroup/metaworld/tree/master/metaworld/policies
[2]https://github.com/sjtuzq/Cycle_Dynamics

**Transferred Policy (Ours)** We encode the context features with $K = 10$ past transitions. The historical state information is postprocessed as state differences between two consecutive states. The historical transition at one step is concatenation of past 10 actions and past 10 postprocessed states. The historical transition data are fed into an MLP with 3 hidden layers with [256, 128, 64] hidden units and Swish activation. The context vector is of dimension 10. The forward dynamics model is an MLP with 4 hidden layers of 200 hidden units and ReLU activation, predicting the state difference between two consecutive states in the future M=10 steps. The learning rate is 0.001 and the batch size is 1024. The objective function is simply $\mathcal{L}_{forw} + \mathcal{L}_{cont}$ (Equation 1 and Equation 2). With these hyper-parameters settings, we train the forward model $F$ and the context model $C$ for 30 epochs, each epoch with 10K steps.

The action translator $H$ first encodes state and action as 128-dimensional vectors respectively. Then, the state embedding and action embedding is concatenated and fed into an MLP with 3 hidden layers of 256 units and ReLU activations. We train the action translator with a decaying learning rate from 3e-4, 5e-5 to 1e-5, and the batch size is also 1024. With these hyper-parameters settings, we train the action translator for 30 epochs, each epoch with 3,000 steps. The objective function is simply $\mathcal{L}_{trans}$ (Equation 3). After each epoch, the performance of the action translator is also evaluated on the target task via averaging the episode rewards in 100 episodes. Finally, the best epoch performance over the 30 epochs is reported.

**Context-conditioned Action Translator** We also demonstrate the performance of policy transfer on more than two tasks as heatmaps in Fig. 5. The heatmaps demonstrate performance gain when comparing our transferred policy against the source policy on the target task. We calculate the improvement in the average episode rewards for every pair of source-target tasks sampled from the training task set. The tasks in the HalfCheetah environment are $\mathcal{T}^{(1)} \cdots \mathcal{T}^{(5)}$ with different armature values, namely {0.1, 0.2, 0.3, 0.4, 0.5}. The tasks in the Ant environment are $\mathcal{T}^{(1)} \cdots \mathcal{T}^{(4)}$ with different leg crippled, namely {0, 1, 2, 3}. As mentioned above, we apply the TD3 algorithm(Fujimoto, Hoof, and Meger 2018) and dense rewards to learn source policies and mediocre policies for each task in training set. Then we collect 100K transition data on each training tasks with the corresponding mediocre policies.

The architecture of context model $C$ and the forward model $F$ remains the same as above, while the learning rate is kept as 5e-4 instead. The architecture of action translator $H$ is expanded to condition on the source task embeddings and target task embeddings. As mentioned in Sec. 3.2, in order to get the representative task feature for any arbitrary training task, we sample 1024 historical transition samples on this task, calculate the their context embedding through context model $C$ and average the 1024 context embedding to get the task feature as an 10-dimensional context vector. The source target feature and target task feature are then encoded as 128-dimensional vectors respectively via a linear layer with ReLU activation. Then the state embedding, action embedding, source task embedding and target task embedding are concatenated to produce the translated action via an MLP with 3 linear layers of 256 hidden units and ReLU activation. The learning rate and batch size for $H$ are 3e-4 and 1024. With these hyper-parameters settings, we train the action translator with 100 epochs, each with 1,000 steps. We report the percentage gain comparing well-trained transferred policies with source policies on each pair of source-target tasks.

## C.3 Policy transfer on tasks sharing a general reward function, differing in dynamics

As explained in Sec. 4, many real-world sparse-reward tasks fall under the umbrella of Proposition 1. Thus, we are mainly interested in policy transfer across tasks with the same reward function $r(s, s')$ but different dynamics. To solve policy transfer across these tasks, our objective function $\mathcal{L}_{trans}$ can be applied so that the transferred policy achieves a value on the target task similar to the source policy on the source task. Experiments in Sec. 5 validate the efficacy of $\mathcal{L}_{trans}$ for learning policy transfer.

As for a more general case, we further consider tasks with different dynamics that have **the same state space, action space and reward function, where the general reward function** $r(s, a, s')$ **cannot be expressed as** $r(s, s')$. Theorem 1 in Appendix A covers this scenario. For source task $\mathcal{T}^{(j)} = \{\mathcal{S}, \mathcal{A}, p^{(j)}, r, \gamma, \rho_0\}$ and target task $\mathcal{T}^{(i)} = \{\mathcal{S}, \mathcal{A}, p^{(i)}, r, \gamma, \rho_0\}$, we can bound the value difference between source policy $\pi^{(j)}$ and transferred policy $\pi^{(i)}$ by minimizing both reward difference $|r(s, \pi^{(i)}(s)) - r(s, \pi^{(j)}(s))|$ and total-variation difference in next state distribution $D_{TV}(p^{(i)}(\cdot|s, \pi^{(i)}(s)), p^{(j)}(\cdot|s, \pi^{(j)}(s))$. Accordingly, we modify transfer loss $\mathcal{L}_{trans}$ (Equation 3) with an additional term of reward difference.

Formally, $\mathcal{L}_{trans,r} = |r_t^{(j)} - R(s_t^{(j)}, \tilde{a}^{(i)}, z^{(i)})| - \lambda \log F(s_{t+1}^{(j)}|s_t^{(j)}, \tilde{a}^{(i)}, z^{(i)})$, where $R$ is a learned reward prediction model, $\lambda$ is a hyper-parameter weight of next state distribution loss, and $\tilde{a}^{(i)} = H(s_t^{(j)}, a_t^{(j)}, z^{(j)}, z^{(i)})$ is the translated action. This objective function drives the action translator H to find an action on the target task leading to a reward and next state, similarly to the source action on the source task.

As explained in Appendix C.1.1, MuJoco environments award the agent considering its velocity of moving forward $v_{torso}$ and the control cost $||a||^2$, i.e. $r = v_{torso} - c||a||^2$. If the coefficient $c = 0$, we can simplify this reward function as $r(s, s')$ because $v_{torso}$ is calculated only based on the current state $s$ and next state $s'$. If $c > 0$, $r$ becomes a general reward function $r(s, a, s')$. We evaluate our action translator trained with $\mathcal{L}_{trans}$ and $\mathcal{L}_{trans,r}$ for this general case of reward function. We search the hyper-parameter value of $\lambda$ in $\mathcal{L}_{trans,r}$ and $\lambda = 10$ performs well across settings.

Our action translator with either $\mathcal{L}_{trans}$ or $\mathcal{L}_{trans,r}$ performs well for policy transfer. When the rewards depend on the action more heavily (i.e. $c$ becomes larger), the advantage of $\mathcal{L}_{trans,r}$ becomes more obvious. However, ours with $\mathcal{L}_{trans,r}$ requires the extra complexity of learning a reward prediction model $R$. When the reward function is mostly determined by the states and can be approximately simplified as $r(s, s')$, we recommend $\mathcal{L}_{trans}$ because it is simpler and achieves a competitive performance.

| Control cost coefficient | Source policy on source task | Source policy on target task | Transferred policy (Zhang et al. 2020c) on target task | Transferred policy (ours with $\mathcal{L}_{trans}$) on target task | Transferred policy (ours with $\mathcal{L}_{trans,r}$) on target task |
|---|---|---|---|---|---|
| c=0.001 | 511.1 | 54.7 | 133.27 | 193.7 | **203.1** |
| c=0.002 | 488.4 | 53.7 | 129.86 | 179.3 | **195.4** |
| c=0.005 | 475.8 | 38.9 | 112.36 | 148.5 | **171.8** |

Table 7: Average episode rewards on Ant environments. We consider the settings with different coefficients for control cost.

On Hopper and HalfCheetah, the control cost coefficient is $c > 0$ by default. Our proposed policy transfer and MCAT achieve performance superior to the baselines on these environments (Sec. 5). This verifies the merits of our objective function $\mathcal{L}_{trans}$ on tasks with a general reward function $r(s, a, s')$.

### C.4   Implementation Details for Comparison with Context-based Meta-RL Algorithms

**Environment**   We modify the physics parameters in the environments to get multiple tasks with varying dynamics functions. We delay the environment rewards to make sparse-reward tasks so that the baseline methods may struggle in these environments. The episode length is set as 1000 steps. The details of the training task set and test task set are shown in Table 8.

| Environment | Reward Delay Steps | Physics Parameter | Training Tasks | Test Tasks |
|---|---|---|---|---|
| Hopper | 100 | Size | {0.02, 0.03, 0.04, 0.05, 0.06} | {0.01, 0.07} |
| HalfCheetah | 500 | Armature | {0.2, 0.3, 0.4, 0.5, 0.6} | {0.05,0.1,0.7,0.75} |
| | | Mass | {0.5, 1.0, 1.5, 2.0, 2.5} | {0.2, 0.3, 2.7, 2.8} |
| Ant | 500 | Damping | {1.0, 10.0, 20.0, 30.0} | {0.5,35.0} |
| | | Crippled Leg | { No crippled leg, crippled leg 0, 1, 2} | {crippled leg 3} |

Table 8: Modified physics parameters used in the experiments.

**Implementation Details**   In Section 5.2, we compare our proposed method with other context-based meta-RL algorithms on environments with sparse rewards. Below we describe the implementation details of each method.

**PEARL(Rakelly et al. 2019)**   We use the implementation provided by the authors[3]. The PEARL agent consists of the context encoder model and the policy model. Following the default setup, the context encoder model is an MLP encoder with 3 hidden layers of 200 units each and ReLU activation. We model the policy as Gaussian, where the mean and log variance is also parameterized by MLP with 3 hidden layers of 300 units and ReLU activation. Same to the default setting, the log variance is clamped to [-2, 20]. We mostly use the default hyper-parameters and search the dimension of the context vector in $\{5, 10, 20\}$. We report the performance of the best hyper-parameter, which achieves highest average score on training tasks.

**MQL(Fakoor et al. 2019)**   We use the implementation provided by the authors[4]. The context encoder is a Gated Recurrent Unit model compressing the information in recent historical transitions. The actor network conditioning on the context features is an MLP with 2 hidden layers of 300 units each and a ReLU activation function. The critic network is of the same architecture as the actor network. We search the hyper-parameters: learning rate in $\{0.0003, 0.0005, 0.001\}$, history length in $\{10, 20\}$, GRU hidden units in $\{20, 30\}$, TD3 policy noise in $\{0.1, 0.2\}$, TD3 exploration noise in $\{0.1, 0.2\}$. We report the performance of the best set of hyper-parameters, which achieves highest average score on training tasks.

**Distral(Teh et al. 2017)**   We use the implementation in the MTRL repository[5]. The Distral framework consists of a central policy and several task-specific policies. The actor network of the central policy is an MLP with 3 hidden layers of 400 units each and a ReLU activation function. The actor and critic networks of the task-specific policies are of the same architecture as the actor network of the central policy. As for the hyperparameters, we set $\alpha$ to 0.5 and search $\beta$ in $\{1, 10, 100\}$, where $\frac{\alpha}{\beta}$ controls the divergence between central policy and task-specific policies, and $\frac{1}{\beta}$ controls the entropy of task-specific policies. When optimizing the actor and critic networks, the learning rates are 1e-3. We report the performance of the best hyper-parameter, which achieves highest average score on training tasks.

---

[3]https://github.com/katerakelly/oyster
[4]https://github.com/amazon-research/meta-q-learning
[5]https://github.com/facebookresearch/mtrl

**HiP-BMDP(Zhang et al. 2020b)**   We use the implementation in the MTRL repository (same as the Distral baseline above). The actor and critic networks are also the same as the ones in Distral above. When optimizing the actor and critic network, the learning rates for both of them are at 1e-3. The log variance of the policy is bound to [-20, 2]. We search the $\Theta$ learning error weight $\alpha_\psi$ in $\{0.01, 0.1, 1\}$, which scales their task bisimulation metric loss. We report the performance of the best hyper-parameter, which achieves highest average score on training tasks.

**MCAT (Ours)**   The architectures of the context model $C$, forward dynamics model $F$ and the action translator $H$ are the same as introduced in Appendix C.2. The actor network and critic network are both MLPs with 2 hidden layers of 256 units and ReLU activations. As described in Algorithm 1, at each iteration, we collect 5K transition data from training tasks. Then we train the context model $C$ and forward dynamics model $F$ for 10K training steps. We train the action translator $H$ for 1K training steps. The actor and critic networks are updated for 5K training steps. In order to monitor the performance of transferred and learned policy in recent episodes, we clear the information about episode reward in $R^{(i)}$ and $R^{(j)\to(i)}$ before the last $G = 20000$ steps.

The learning rate and batch size of training $C$, $F$ and $H$ are the same as introduced in "Context-conditioned Action Translator" in Appendix C.2. The hyper-parameters of learning the actor and critic are the same as listed in Table 5. Besides, we adapt the official implementation[6] to maintain SIL replay buffer with their default hyper-parameters on MuJoCo environments.

Even though there are a couple of components, they are trained alternatively not jointly. The dynamics model is learned with $\mathcal{L}_{forw}$ to accurately predict the next state. The learned context embeddings for different tasks can separate well due to the regularization term $\mathcal{L}_{const}$. With the fixed context encoder and dynamics model, the action translator can be optimized. Then, with the fixed context encoder, the context-conditioned policy learns good behavior from data collected by the transferred policy. These components are not moving simultaneously and this fact facilitates the learning process. To run our approach on MuJoCo environments, for each job, we need to use one GPU card (NVIDIA GeForce GTX TITAN X) for around 4 days. Fig. 6 and Tab. 9 show the performance of our approach and baselines on various environments.

| Setting | Hopper Size | HalfCheetah Armature | HalfCheetah Mass | Ant Damping | Ant Cripple |
|---|---|---|---|---|---|
| MQL | 1607.5 ($\pm$327.5) | -77.9 ($\pm$214.2) | -413.9 ($\pm$11.0) | 103.1 ($\pm$35.7) | 38.2 ($\pm$4.0) |
| PEARL | 1755.8 ($\pm$115.3) | -18.8 ($\pm$69.3) | 25.9 ($\pm$69.2) | 73.2 ($\pm$13.3) | 3.5 ($\pm$2.4) |
| Distral | 1319.8 ($\pm$162.2) | 566.9 ($\pm$246.7) | -29.5 ($\pm$3.0) | 90.5 ($\pm$28.4) | -0.1 ($\pm$0.7) |
| HiP-BMDP | 1368.3 ($\pm$150.7) | -102.4 ($\pm$24.9) | -74.8 ($\pm$35.4) | 33.1 ($\pm$6.0) | 7.3 ($\pm$2.6) |
| MCAT(Ours) | **1914.8** ($\pm$373.2) | **2071.5** ($\pm$447.4) | **1771.1** ($\pm$617.7) | **624.6** ($\pm$218.8) | **281.6** ($\pm$65.6) |

Table 9: Test rewards at 2M timesteps, averaged over 3 runs. This corrects the last column of Tab. 2 in the main text.
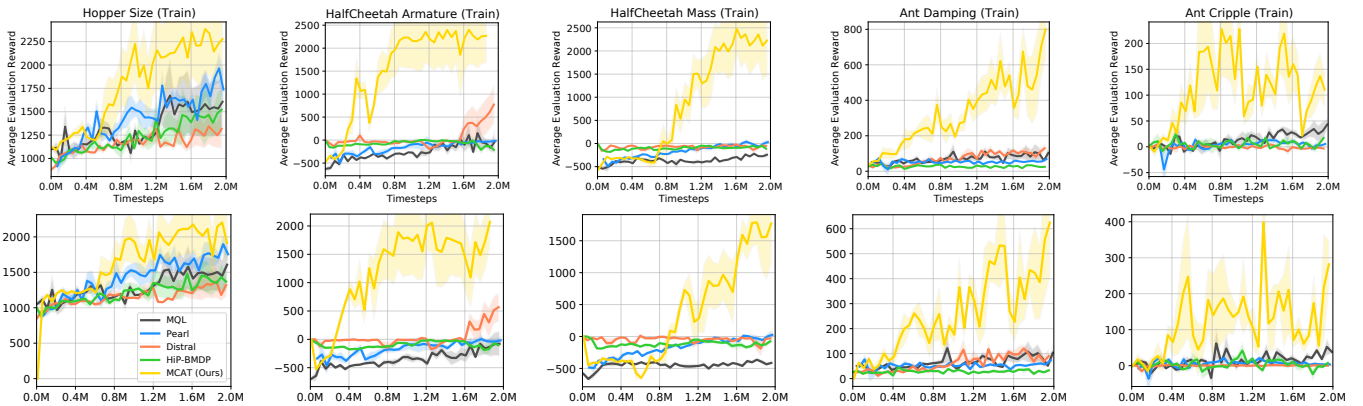


Figure 6: Learning curves of episode rewards on both training and test tasks, averaged over 3 runs. Shadow areas indicate standard error. This adds the performance on training tasks in comparison to Fig. 3 in the main text and changes the learning curves on test tasks in Ant Cripple setting to 2M timesteps

---

[6]https://github.com/junhyukoh/self-imitation-learning

Furthermore, we present additional experimental results on MetaWorld environment. In Section 5.1, we introduced the tasks of moving objects to target locations and the reward is positive only when the object is close to the goal. We combine context-based TD3 with policy transfer to learn a policy operating multiple objects: drawer, coffee mug, soccer, cube, plate. Then we test whether the policy could generalize to moving a large cylinder. In Tab. 10, MCAT agent earns higher success rate than the baselines on both training and test tasks after 2M timesteps in the sparse-reward tasks.

|  | MQL (Fakoor et al. 2019) | PEARL (Rakelly et al. 2019) | PCGrad (Yu et al. 2020a) | MCAT |
| --- | --- | --- | --- | --- |
| Training tasks (reward) | $164.8_{(\pm23.6)}$ | $161.2_{(\pm25.3)}$ | $44.8_{(\pm31.7)}$ | $\mathbf{204.1}_{(\pm43.1)}$ |
| Test tasks (reward) | $0.0_{(\pm0.0)}$ | $0.0_{(\pm0.0)}$ | $0.0_{(\pm0.0)}$ | $\mathbf{10.2}_{(\pm8.3)}$ |
| Training tasks (success rate) | $40.0\%_{(\pm0.0\%)}$ | $33.3\%_{(\pm5.4\%)}$ | $10.0\%_{(\pm7.1\%)}$ | $\mathbf{53.3\%}_{(\pm5.4\%)}$ |
| Test tasks (success rate) | $0.0\%_{(\pm0.0\%)}$ | $0.0\%_{(\pm0.0\%)}$ | $0.0\%_{(\pm0.0\%)}$ | $\mathbf{16.7\%}_{(\pm13.6\%)}$ |

Table 10: Performance of learned policies at 2M timesteps, averaged over 3 runs.

# D Ablative Study

## D.1 Effect of Policy Transfer

In Section 5.3, we investigate the effect of policy transfer (PT). In Figure 7 we provide the learning curves of MCAT and MCAT w/o PT on both training tasks and test tasks.
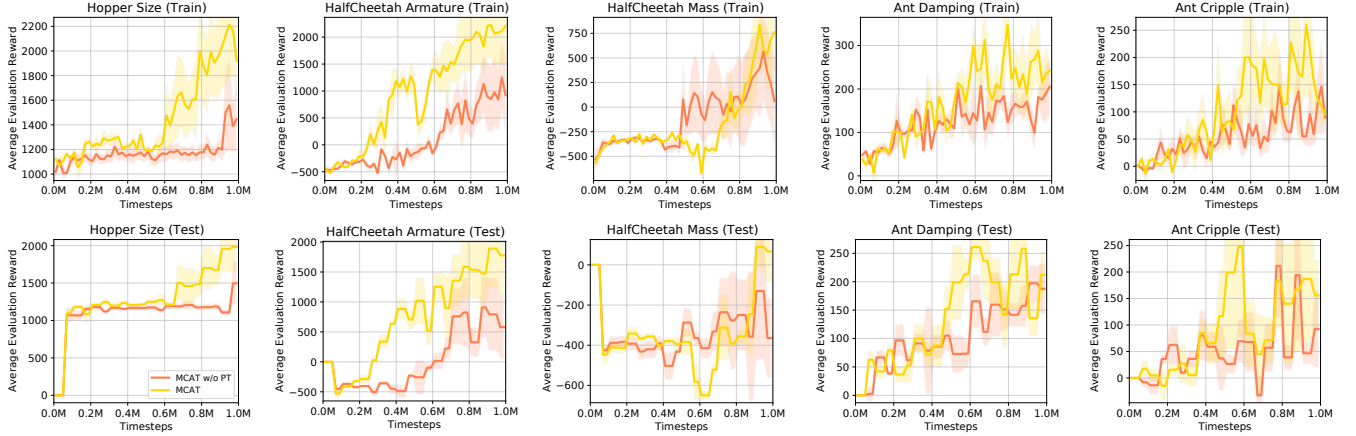


Figure 7: Learning curves of the average episode reward, averaged over 3 runs. The average episode reward and standard error are reported on training tasks and test tasks respectively. This repeats Figure 7 with addition of learning curves on training tasks.

## D.2 More Sparse Rewards

In Section 5.3, we report the effect of policy transfer when the rewards become more sparse in the environments. On HalfChee-tah, we delay the environment rewards for different number of steps 200, 350, 500. In Figure 8, we show the learning curves on training and test tasks. In Table 4, we report the average episode rewards and standard error over 3 runs at 1M timesteps.
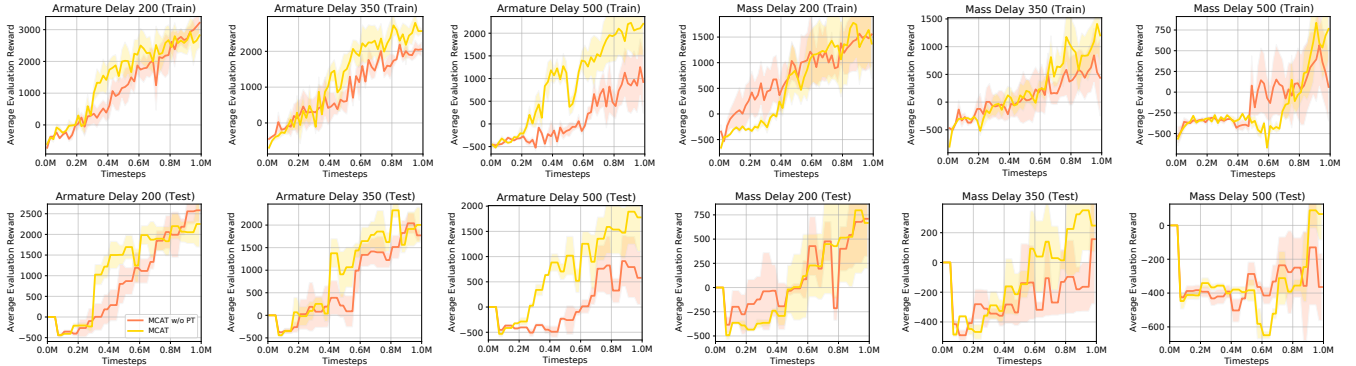


Figure 8: Learning curves of the average episode reward, averaged over 3 runs. The average episode reward and standard error are reported on training tasks and test tasks respectively.

## D.3 More Diverse Tasks

We include more settings of training and test tasks where the discrepancy among training tasks varies. On HalfCheetah, the environment rewards are delayed for 500 steps. In Table 11, we list the details of the settings.

| Physics Parameter | Setting | Train | Test |
|---|---|---|---|
| Armature | Set 1 | {0.2, 0.25, 0.3, 0.35, 0.4} | {0.05, 0.1, 0.5, 0.55} |
|  | Set 2 | {0.2, 0.3, 0.4, 0.5, 0.6} | {0.2, 0.3, 0.7, 0.75} |
|  | Set 3 | {0.2, 0.35, 0.5, 0.65, 0.8} | {0.2, 0.3, 0.9, 0.95} |
| Mass | Set 1 | {0.5, 0.75, 1.0, 1.25, 1.5} | {0.2, 0.3, 1.7, 1.8} |
|  | Set 2 | {0.5, 1.0, 1.5, 2.0, 2.5} | {0.2, 0.3, 2.7, 2.8} |
|  | Set 3 | {0.5, 1.25, 2.0, 2.75, 3.5} | {0.2, 0.3, 3.7, 3.8} |

Table 11: Modified physics parameters used in the experiments.

We consider baseline MQL because it performs reasonably well on HalfCheetah among all the baselines (Figure 3). Table 12 demonstrates that policy transfer (PT) is generally and consistently effective. In Figure 9, we show the learning curves on training and test tasks. In Table 12, we report the average episode rewards and standard error over 3 runs at 1M timesteps.

| Setting | Armature Set 1 | | Armature Set 2 | | Armature Set 3 | | Mass Set 1 | | Mass Set 2 | | Mass Set 3 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Task | Train | Test | Train | Test | Train | Test | Train | Test | Train | Test | Train | Test |
| MQL | -129.3 (±46.7) | -248.0 (±32.0) | -277.2 (±25.2) | -335.0 (±20.8) | -85.0 (±33.5) | -214.7 (±28.9) | -100.8 (±37.8) | -291.3 (±25.8) | -403.7 (±16.1) | -467.8 (±6.5) | -175.3 (±6.2) | -287.9 (±11.7) |
| MCAT w/o PT | 837.6 (±646.5) | 785.3 (±733.1) | 924.0 (±690.1) | 579.1 (±527.1) | 452.8 (±386.6) | 616.5 (±305.0) | -60.5 (±313.4) | -258.2 (±151.1) | 62.5 (±411.0) | -364.3 (±198.5) | -328.1 (±55.8) | -412.4 (±7.7) |
| MCAT | 3372.1 (±186.4) | 2821.9 (±137.7) | 2207.3 (±697.7) | 1776.8 (±680.8) | 1622.2 (±402.2) | 918.3 (±142.5) | 1222.2 (±754.9) | 482.4 (±624.2) | 763.4 (±377.7) | 67.1 (±152.9) | 705.7 (±503.4) | -86.2 (±111.8) |
| Improvement(%) | 302.6 | 259.3 | 133.9 | 206.8 | 258.3 | 49.0 | 2120.2 | 286.8 | 1121.4 | 118.4 | 315.1 | 79.1 |

Table 12: The performance of learned policy on various task settings. We modify *armature* and *mass* to get 5 training tasks and 4 test tasks in each setting. We compute the improvement of MCAT over MCAT w/o PT.
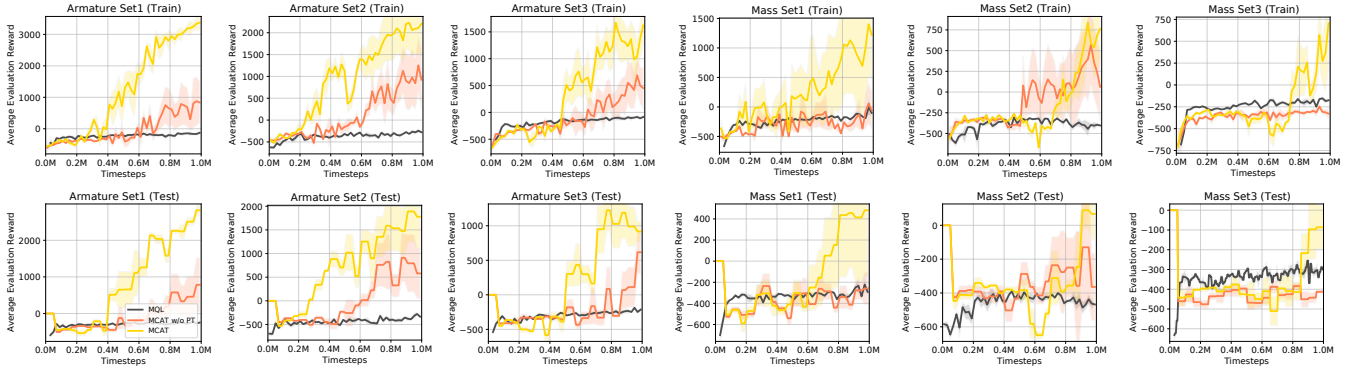


Figure 9: Learning curves of the average episode reward, averaged over 3 runs. The average episode reward and standard error are reported on training tasks and test tasks respectively.

## D.4 Effect of Self-Imitation Learning

We run experiments combining baseline methods with self-imitation learning (SIL) (Oh et al. 2018). SIL brings improvement to baselines but still ours shows significant advantages. In Tab. 13, MCAT w/o SIL compares favorably with the baseline methods. MCAT further improves the performance of MCAT w/o SIL, and MCAT outperform the variants of baseline methods with SIL.

| Setting | Hopper Size | | HalfCheetah Armature | | HalfCheetah Mass | | Ant Damping | | Ant Cripple | |
|---|---|---|---|---|---|---|---|---|---|---|
| Task | Train | Test | Train | Test | Train | Test | Train | Test | Train | Test |
| MQL (Fakoor et al. 2019) | 1586.1 (± 321.4) | 1607.5 (± 327.5) | -31.4 (± 243.5) | -77.9 (± 214.3) | -243.1 (± 69.8) | -413.9 (± 11.1) | 93.8 (± 24.5) | 103.1 (± 35.7) | 17.4 (± 4.3) | 38.2 (± 4.0) |
| Distral (Teh et al. 2017) | 1364.0 (± 216.3) | 1319.8 (± 162.2) | 774.7 (± 405.9) | 566.9 (± 246.7) | -54.3 (± 14.8) | -29.5 (± 3.0) | 123.0 (± 20.0) | 90.5 (± 28.4) | -2.5 (± 1.7) | -0.1 (± 0.7) |
| HiP-BMDP (Zhang et al. 2020b) | 1590.3 (± 238.7) | 1368.3 (± 150.7) | -212.4 (± 52.2) | -102.4 (± 24.9) | -81.3 (± 8.31) | -101.8 (± 29.6) | 15.0 (± 5.7) | 33.1 (± 6.0) | 12.7 (± 5.3) | 7.3 (± 2.6) |
| MCAT w/o SIL | 1261.6 (± 55.2) | 1165.1 (± 8.6) | 1548.8 (± 418.4) | 883.8 (± 267.2) | 610.6 (± 482.3) | 119.0 (± 210.0) | 123.3 (± 25.8) | 123.8 (± 26.9) | 97.3 (± 3.6) | 163.1 (± 26.1) |
| MQL+SIL | 1395.5 (± 60.8) | 1398.9 (± 85.9) | 1399.7 (± 350.2) | 743.5 (± 246.1) | 617.8 (± 133.1) | -63.3 (± 158.3) | 153.0 (± 28.3) | 144.3 (± 28.1) | 13.9 (± 19.8) | 10.2 (± 2.3) |
| Distral+SIL | 1090.2 (± 18.7) | 1090.9 (± 7.8) | 1014.1 (± 121.4) | 970.3 (± 164.2) | 809.7 (± 294.2) | 746.7 (± 120.5) | 174.3 (± 66.1) | 122.2 (± 44.5) | 107.7 (± 57.7) | 9.1 (± 5.0) |
| HiP-BMDP+SIL | 1573.3 (± 32.4) | 1589.5 (± 110.3) | 954.8 (± 192.3) | 713.3 (± 85.4) | 953.5 (± 61.2) | 506.6 (± 99.0) | 653.9 (± 262.6) | 523.6 (± 300.8) | **170.9** (± 68.7) | 215.4 (± 130.3) |
| MCAT (Ours) | **2278.8** (± 426.2) | **1914.8** (± 373.2) | **2267.2** (± 579.2) | **2071.5** (± 447.4) | **2226.3** (± 762.6) | **1771.1** (± 617.7) | **1322.7** (± 57.4) | **1014.0** (± 69.9) | 110.4 (± 30.5) | **281.6** (± 65.6) |

Table 13: Mean (± standard error) of episode rewards on the training and test tasks, at 2M timesteps.

On one task, SIL boosts the performance by exploiting the successful past experiences. But on multiple tasks, enhancing performance on one task with luckily collected good experiences may not benefit the exploration on other tasks. If other tasks have never seen the good performance before, SIL might even prevent the exploration on these tasks because the shared policy is trained to overfit highly-rewarding transitions on the one task with good past trajectories. We observe that after combining with SIL, the baselines show even more severe performance imbalance among multiple training tasks. Therefore, we believe the idea of policy transfer is complementary to SIL, in that it makes each task benefit from good policies on any other tasks.

## D.5 Effect of Contrastive Loss

To show the contrastive loss indeed helps policy transfer, we compare our method with and without the contrastive loss $\mathcal{L}_{cont}$ ( Equation 2). In Fig. 10, one can observe that $\mathcal{L}_{cont}$ helps cluster the embeddings of samples from the same task and separate the embeddings from different tasks. We note that the tasks $\mathcal{T}^{(1)}, \mathcal{T}^{(2)}, \mathcal{T}^{(3)}, \mathcal{T}^{(4)}, \mathcal{T}^{(5)}$ have different values of the physics parameter armature $0.2, 0.3, 0.4, 0.5, 0.6$. As mentioned in Sec. 3.2, the learned context embeddings maintain the similarity between tasks. In Fig. 10, the context embeddings of two tasks are closer if their values of armature is closer.



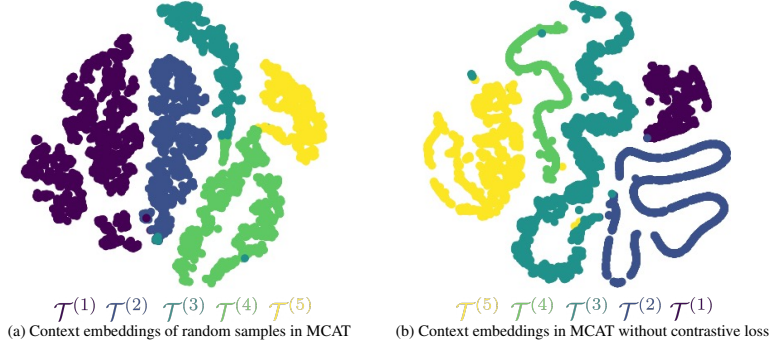(a) Context embeddings of random samples in MCAT    (b) Context embeddings in MCAT without contrastive loss

Figure 10: t-SNE visualization(Van der Maaten and Hinton 2008) of the context embeddings learned via our method with and without contrastive loss. Different colors correspond to different training tasks.

MCAT shows superior performance to the variant without the contrastive loss. Here we show the learning curves on training and test tasks separately(Fig. 11).
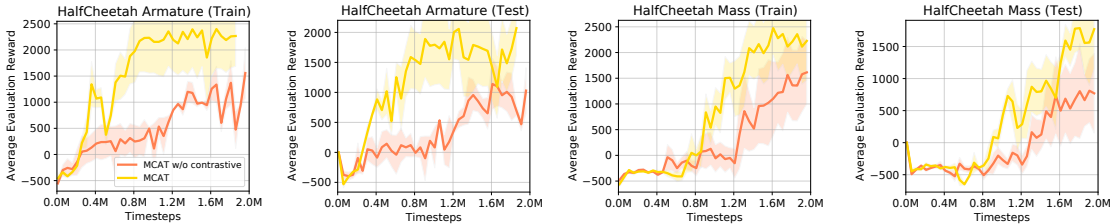


Figure 11: Learning curves of the average episode reward, averaged over 3 runs. The average episode reward and standard error are reported on training tasks and test tasks respectively.

## D.6 Design Choice of Action Translator

We add this experimental comparison with the action translator by (Zhang et al. 2020c). To learn a shared policy solving multiple tasks, we combine the context-based TD3 algorithm, self-imitation learning, and policy transfer with their action translator. Using their action translator underperforms ours. The main reason is that, with changing datasets and policies, their action translator may be harder to tune because there are more moving components (i.e. another action translator, a discriminator) and more loss terms to be balanced (i.e. domain cycle-consistency loss, adversarial loss).

| Setting | HalfCheetah Armature | | HalfCheetah Mass | |
|---|---|---|---|---|
| Tasks | Training | Test | Training | Test |
| MCAT | **2267.2** $(\pm 579.2)$ | **2071.5** $(\pm 447.4)$ | **2226.3** $(\pm 762.6)$ | **1771.1** $(\pm 617.7)$ |
| MCAT with (Zhang et al. 2020c) action translator | 2255.2 $(\pm 644.4)$ | 1664.8 $(\pm 660.8)$ | 1185.8 $(\pm 798.0)$ | 684.7 $(\pm 759.0)$ |

Table 14: Mean ($\pm$ standard error) of episode rewards on training and test tasks at 2M timesteps.

# E  Extension of Policy Transfer

As clarified in Sec. 3.1, in this work, we mainly focus on tasks with the same state space, action space, reward function but varying dynamics. However, we note that our proposed method of learning action translator may be extended to tackle the challenge of policy transfer in more general cases, such as (1) Tasks differing in reward function, (2) Tasks differing in state space and action space. In this section, we establish the theory and method in details to extend our policy transfer approach, as a supplement to Sec. 6.

## E.1  Theoretical Analysis

Intuitively, on two general tasks, we aim to discover correspondent state-action pairs achieving the same reward and transiting to correspondent next states. With the state and action correspondence, the behavior of good source policy can be "replicated" in the target task and the high value of the good source policy can be maintained by the transferred policy on the target task. Inspired by this idea, we extend our theory in Sec. 4 and Appendix A.

We first define a binary relation for states to describe the equivalent states on two MDPs (Definition 2) and define an invertible function to capture the state equivalence relation (Definition 3). Building upon the notion of state equivalence, we derive the upper bound of value difference between policies on two MDPs (Theorem 2). Finally, we reach a proposition for the upper bound of value difference (Proposition 2) to explain that our objective function in learning action translator can be extended to bound the value difference between the source and transferred policy.

**Definition 2.** *Given two MDPs $\mathcal{T}^{(i)} = \{\mathcal{S}^{(i)}, \mathcal{A}^{(i)}, p^{(i)}, r^{(i)}, \gamma, \rho_0^{(i)}\}$ and $\mathcal{T}^{(j)} = \{\mathcal{S}^{(j)}, \mathcal{A}^{(j)}, p^{(j)}, r^{(j)}, \gamma, \rho_0^{(j)}\}$, we define a binary relation $B \in \mathcal{S}^{(i)} \times \mathcal{S}^{(j)}$ called **state equivalence relation**. Let $s'^{(i)}$ denote the next state following state $s^{(i)}$, and $s'^{(j)}$ denote the next state following state $s^{(j)}$. For states $s^{(i)} \in \mathcal{S}^{(i)}$, $s^{(j)} \in \mathcal{S}^{(j)}$, we have $(s^{(i)}, s^{(j)}) \in B$ (i.e. $s^{(i)} B s^{(j)}$) if for any $a^{(i)} \in \mathcal{A}^{(i)}$ there exists $a^{(j)} \in \mathcal{A}^{(j)}$ satisfying the following conditions:*

$$r^{(i)}(s^{(i)}, a^{(i)}) = r^{(j)}(s^{(j)}, a^{(j)})$$

$$\forall s'^{(i)} \in S^{(i)}, \exists s'^{(j)} \in \mathcal{S}^{(j)} \text{ s.t. } p^{(i)}(s'^{(i)}|s^{(i)}, a^{(i)}) = p^{(j)}(s'^{(j)}|s^{(i)}, a^{(j)}) \text{ and } s'^{(i)} B s'^{(j)}$$

We call the state $s^{(i)}$ and $s^{(j)}$ are correspondent/equivalent when $(s^{(i)}, s^{(j)}) \in B$. Also, in this case, the action $a^{(i)}$ for state $s^{(i)}$ on the MDP $\mathcal{T}^{(i)}$ is equivalent to the action $a^{(j)}$ for state $s^{(j)}$ on the MDP $\mathcal{T}^{(j)}$.

This definition is related to stochastic bisimulation relation in (Ferns, Panangaden, and Precup 2004; Zhang et al. 2020a,b). Unlike these prior works about state bisimulation, we allow the equivalent actions $a^{(j)} \neq a^{(i)}$. So action $a$ on the task $\mathcal{T}^{(i)}$ might not be equivalent to $a$ on the task $\mathcal{T}^{(j)}$, and hence we need to involve action translator in learning of both the state correspondence and action correspondence.

Drawing upon Definition 2, we define a one-to-one mapping to identify the equivalent state across two spaces $\mathcal{S}^{(i)}$ and $\mathcal{S}^{(j)}$.

**Definition 3.** *Given two MDPs $\mathcal{T}^{(i)} = \{\mathcal{S}^{(i)}, \mathcal{A}^{(i)}, p^{(i)}, r^{(i)}, \gamma, \rho_0^{(i)}\}$ and $\mathcal{T}^{(j)} = \{\mathcal{S}^{(j)}, \mathcal{A}^{(j)}, p^{(j)}, r^{(j)}, \gamma, \rho_0^{(j)}\}$ with state equivalence relation $B$, we consider subsets $\mathcal{S}_B^{(i)} \subset \mathcal{S}^{(i)}$ and $\mathcal{S}_B^{(j)} \subset \mathcal{S}^{(j)}$ satisfying: $\forall s^{(i)} \in \mathcal{S}_B^{(i)}, \exists s^{(j)} \in \mathcal{S}_B^{(j)}$ s.t. $(s^{(i)}, s^{(j)}) \in B$. We define a invertible function $G : \mathcal{S}_B^{(i)} \to \mathcal{S}_B^{(j)}$ called **state translator function**, satisfying: $(s^{(i)}, G(s^{(i)})) \in B$.*

Based on Defintion 2 and 3, given two correspondent states $s^{(i)} \in S_B^{(i)}$ and $s^{(j)} \in S_B^{(j)}$, we can derive the upper bound for the value difference between $V^{\pi^{(i)}}(s^{(i)}, \mathcal{T}^{(i)})$ and $V^{\pi^{(j)}}(s^{(j)}, \mathcal{T}^{(j)})$.

**Theorem 2.** *$\mathcal{T}^{(i)} = \{\mathcal{S}^{(i)}, \mathcal{A}^{(j)}, p^{(i)}, r^{(i)}, \gamma, \rho_0^{(i)}\}$ and $\mathcal{T}^{(j)} = \{\mathcal{S}^{(j)}, \mathcal{A}^{(j)}, p^{(j)}, r^{(j)}, \gamma, \rho_0^{(j)}\}$ are two MDPs sampled from the distribution of tasks $p(\mathcal{T})$. $\pi^{(i)}$ is a deterministic policy on $\mathcal{T}^{(i)}$ and $\pi^{(j)}$ is a deterministic policy on $\mathcal{T}^{(j)}$. We assume there exist state equivalence relation $B \in \mathcal{S}^{(i)} \times \mathcal{S}^{(j)}$ and a state translator function $G$ defining a one-to-one mapping from $\mathcal{S}_B^{(i)}$ to $\mathcal{S}_B^{(j)}$. Let $M = \sup_{s^{(i)} \in \mathcal{S}^{(i)}} |V^{\pi^{(i)}}(s^{(i)}, \mathcal{T}^{(i)})|$ and $d = \sup_{s^{(i)} \in \mathcal{S}_B^{(i)}} \left[|r^{(i)}(s^{(i)}, \pi^{(i)}(s)) - r^{(j)}(s^{(j)}, \pi^{(j)}(s))| + 2\gamma M D_{TV}(p^{(i)}(\cdot|s^{(i)}, \pi^{(i)}(s^{(i)})), p^{(j)}(G(\cdot)|s^{(j)}, \pi^{(j)}(s^{(j)})))\right]$. Then $\forall s^{(i)} \in \mathcal{S}_B^{(i)}, s^{(j)} = G(s^{(i)})$, we have*

$$\left| V^{\pi^{(i)}}(s^{(i)}, \mathcal{T}^{(i)}) - V^{\pi^{(j)}}(s^{(j)}, \mathcal{T}^{(j)}) \right| \leq \frac{d}{1 - \gamma}$$

*Proof.* Let $a^{(i)} = \pi^{(i)}(s^{(i)})$ and $a^{(j)} = \pi^{(j)}(s^{(j)})$. We rewrite the value difference.

$$V^{\pi^{(i)}}(s^{(i)}, \mathcal{T}^{(i)}) - V^{\pi^{(j)}}(s^{(j)}, \mathcal{T}^{(j)})$$

$$= r^{(i)}(s^{(i)}, a^{(i)}) + \gamma \sum_{s'^{(i)} \in \mathcal{S}^{(i)}} p^{(i)}(s'^{(i)}|s^{(i)}, a^{(i)})V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)}) - r^{(j)}(s^{(j)}, a^{(j)}) - \gamma \sum_{s'^{(j)} \in \mathcal{S}^{(j)}} p^{(j)}(s'^{(j)}|s^{(j)}, a^{(j)})V^{\pi^{(j)}}(s'^{(j)}, \mathcal{T}^{(j)})$$

$$= (r^{(i)}(s^{(i)}, a^{(i)}) - r^{(j)}(s^{(j)}, a^{(j)}))$$

$$+ \gamma(\sum_{s'^{(i)} \in \mathcal{S}^{(i)}} p^{(i)}(s'^{(i)}|s^{(i)}, a^{(i)})V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)}) - \sum_{s'^{(j)} \in \mathcal{S}^{(j)}} p^{(j)}(s'^{(j)}|s^{(j)}, a^{(j)})V^{\pi^{(j)}}(s'^{(j)}, \mathcal{T}^{(j)}))$$

According to Definition 2, since $s^{(i)} \in \mathcal{S}_B^{(i)}$, we have $s'^{(i)} \in \mathcal{S}_B^{(i)}$. Similarly, $s'^{(j)} \in \mathcal{S}_B^{(j)}$.
Then we derive the second term in the right side of the equation above:

$$\sum_{s'^{(i)} \in \mathcal{S}^{(i)}} p^{(i)}(s'^{(i)}|s^{(i)}, a^{(i)})V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)}) - \sum_{s'^{(j)} \in \mathcal{S}^{(j)}} p^{(j)}(s'^{(j)}|s^{(j)}, a^{(j)})V^{\pi^{(j)}}(s'^{(j)}, \mathcal{T}^{(j)})$$

*replace $\mathcal{S}^{(i)}$ by $\mathcal{S}_B^{(i)}$ because $s'^{(i)} \in \mathcal{S}_B^{(i)}$, replace $\mathcal{S}^{(j)}$ by $\mathcal{S}_B^{(j)}$ because $s'^{(j)} \in \mathcal{S}_B^{(j)}$

**minus and plus $\sum_{s'^{(i)} \in \mathcal{S}_B^{(i)}} p^{(j)}(G(s'^{(i)})|s^{(j)}, a^{(j)}))V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)})$

$$= \sum_{s'^{(i)} \in \mathcal{S}_B^{(i)}} p^{(i)}(s'^{(i)}|s^{(i)}, a^{(i)})V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)}) - \sum_{s'^{(i)} \in \mathcal{S}_B^{(i)}} p^{(j)}(G(s'^{(i)})|s^{(j)}, a^{(j)}))V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)})$$

$$+ \sum_{s'^{(i)} \in \mathcal{S}_B^{(i)}} p^{(j)}(G(s'^{(i)})|s^{(j)}, a^{(j)}))V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)}) - \sum_{s'^{(j)} \in \mathcal{S}_B^{(j)}} p^{(j)}(s'^{(j)}|s^{(j)}, a^{(j)})V^{\pi^{(j)}}(s'^{(j)}, \mathcal{T}^{(j)})$$

*combine the first two terms, rewrite the third term because $G$ is invertible function

$$= \sum_{s'^{(i)} \in \mathcal{S}_B^{(i)}} \left[ p^{(i)}(s'^{(i)}|s^{(i)}, a^{(i)}) - p^{(j)}(G(s'^{(i)})|s^{(j)}, a^{(j)}) \right] V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)})$$

$$+ \sum_{s'^{(j)} \in \mathcal{S}_B^{(j)}} p^{(j)}(s'^{(j)}|s^{(j)}, a^{(j)}))V^{\pi^{(i)}}(G^{-1}(s'^{(j)}), \mathcal{T}^{(i)}) - \sum_{s'^{(j)} \in \mathcal{S}_B^{(j)}} p^{(j)}(s'^{(j)}|s^{(j)}, a^{(j)})V^{\pi^{(j)}}(s'^{(j)}, \mathcal{T}^{(j)})$$

*combine the last two terms

$$= \sum_{s'^{(i)} \in \mathcal{S}_B^{(i)}} \left[ p^{(i)}(s'^{(i)}|s^{(i)}, a^{(i)}) - p^{(j)}(G(s'^{(i)})|s^{(j)}, a^{(j)}) \right] V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)})$$

$$+ \sum_{s'^{(j)} \in \mathcal{S}_B^{(j)}} p^{(j)}(s'^{(j)}|s^{(j)}, a^{(j)})) \left[ V^{\pi^{(i)}}(G^{-1}(s'^{(j)}), \mathcal{T}^{(i)}) - V^{\pi^{(j)}}(s'^{(j)}, \mathcal{T}^{(j)}) \right]$$

Therefore, we can bound the absolute value of the value difference according to the two equation arrays above:

$$\left| V^{\pi^{(i)}}(s^{(i)}, \mathcal{T}^{(i)}) - V^{\pi^{(j)}}(s^{(j)}, \mathcal{T}^{(j)}) \right| \leq \left| r^{(i)}(s^{(i)}, a^{(i)}) - r^{(j)}(s^{(j)}, a^{(j)}) \right|$$

$$+ \gamma \sum_{s'^{(i)} \in \mathcal{S}_B^{(i)}} \left[ p^{(i)}(s'^{(i)}|s^{(i)}, a^{(i)}) - p^{(j)}(G(s'^{(i)})|s^{(j)}, a^{(j)}) \right] V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)})$$

$$+ \gamma \sum_{s'^{(j)} \in \mathcal{S}_B^{(j)}} p^{(j)}(s'^{(j)}|s^{(j)}, a^{(j)})) \left[ V^{\pi^{(i)}}(G^{-1}(s'^{(j)}), \mathcal{T}^{(i)}) - V^{\pi^{(j)}}(s'^{(j)}, \mathcal{T}^{(j)}) \right]$$

$$\leq \left| r^{(i)}(s^{(i)}, a^{(i)}) - r^{(j)}(s^{(j)}, a^{(j)}) \right| + 2\gamma M D_{TV}(p^{(i)}(\cdot|s^{(i)}, a^{(i)}), p^{(j)}(G(\cdot)|s^{(j)}, a^{(j)}))$$

$$+ \gamma \sup_{s'^{(j)} \in \mathcal{S}_B^{(j)}} \left| V^{\pi^{(i)}}(G^{-1}(s'^{(j)}), \mathcal{T}^{(i)}) - V^{\pi^{(j)}}(s'^{(j)}, \mathcal{T}^{(j)}) \right|$$

$$\leq d + \gamma \sup_{s'^{(i)} \in \mathcal{S}_B^{(i)}} \left| V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)}) - V^{\pi^{(j)}}(G(s'^{(i)}), \mathcal{T}^{(j)}) \right| \leq \frac{d}{1-\gamma}$$

$\square$

Theorem 2 proves the value difference is upper bounded by a scalar $d$, depending on the reward difference $|r^{(i)}(s^{(i)}, \pi^{(i)}(s^{(i)})) - r^{(j)}(s^{(j)}, \pi^{(j)}(s^{(j)}))|$ and $D_{TV}(p^{(i)}(\cdot|s^{(i)}, a^{(i)}), p^{(j)}(G(\cdot)|s^{(j)}, a^{(j)}))$, i.e. the total-variation distance between probability distribution of next state on $\mathcal{T}^{(i)}$ and probability distribution of correspondent next state on $\mathcal{T}^{(j)}$. Indeed, if the state equivalence relation is only true for identical states (i.e. $G$ is an identity mapping, $s^{(i)} B s^{(j)}$ if and only if $s^{(i)} = s^{(j)}$), then Theorem 2 degenerates into Theorem 1. We note the proof of Theorem 2 is similar to proof of Theorem 1 in Appendix A.

For a special case, where the reward only depends on the current state and next state, we can formulate a simpler definition of scalar $d$. The following Proposition 2 is analogous to Proposition 1 in the assumption about reward function.

**Proposition 2.** $\mathcal{T}^{(i)} = \{\mathcal{S}^{(i)}, \mathcal{A}(i), p^{(i)}, r^{(i)}, \gamma, \rho_0^{(i)}\}$ and $\mathcal{T}^{(j)} = \{\mathcal{S}^{(j)}, \mathcal{A}^{(j)}, p^{(j)}, r^{(j)}, \gamma, \rho_0^{(j)}\}$ *are two MDPs sampled from the distribution of tasks $p(\mathcal{T})$. $\pi^{(i)}$ is a deterministic policy on $\mathcal{T}^{(i)}$ and $\pi^{(j)}$ is a deterministic policy on $\mathcal{T}^{(j)}$. We assume there exist state equivalence relation $B \in \mathcal{S}^{(i)} \times \mathcal{S}^{(j)}$ and a state translator function $G$ defining a one-to-one mapping from $\mathcal{S}_B^{(i)}$ to $\mathcal{S}_B^{(j)}$. Suppose that the reward function $r^{(i)}(s^{(i)}, a^{(i)}, s'^{(i)}) = r^{(i)}(s^{(i)}, s'^{(i)})$ and $r^{(j)}(s^{(j)}, a^{(j)}, s'^{(j)}) = r^{(j)}(s^{(j)}, s'^{(j)})$. If $s^{(j)} = G(s^{(i)})$ and $s'^{(j)} = G(s'^{(i)})$, $r^{(i)}(s^{(i)}, s'^{(i)}) = r^{(j)}(s^{(j)}, s'^{(j)})$. Let $M = \sup_{s^{(i)} \in \mathcal{S}^{(i)}} |r^{(i)}(s^{(i)}, s'^{(i)}) + \gamma V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)})|$ and $d = \sup_{s^{(i)} \in \mathcal{S}_B^{(i)}} 2 M D_{TV}(p^{(i)}(\cdot|s^{(i)}, \pi^{(i)}(s^{(i)})), p^{(j)}(G(\cdot)|s^{(j)}, \pi^{(j)}(s^{(j)})))$.*

*Then $\forall s^{(i)} \in \mathcal{S}_B^{(i)}, s^{(j)} = G(s^{(i)})$, we have*

$$\left| V^{\pi^{(i)}}(s^{(i)}, \mathcal{T}^{(i)}) - V^{\pi^{(j)}}(s^{(j)}, \mathcal{T}^{(j)}) \right| \leq \frac{d}{1 - \gamma}$$

*Proof.* Let $a^{(i)} = \pi^{(i)}(s^{(i)})$ and $a^{(j)} = \pi^{(j)}(s^{(j)})$. $s'^{(i)}$ denotes the next state following state $s^{(i)}$.

Because the reward solely depends on the current and next state, we rewrite the value function:

$$
\begin{aligned}
V^{\pi^{(i)}}(s^{(i)}, \mathcal{T}^{(i)}) &= r^{(i)}(s^{(i)}, a^{(i)}) + \gamma \sum_{s'^{(i)} \in \mathcal{S}^{(i)}} p^{(i)}(s'^{(i)}|s^{(i)}, a^{(i)}) V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)}) \\
&= \sum_{s'^{(i)} \in \mathcal{S}^{(i)}} p^{(i)}(s'^{(i)}|s^{(i)}, a^{(i)}) r^{(i)}(s^{(i)}, s'^{(i)}) + \gamma \sum_{s'^{(i)} \in \mathcal{S}^{(i)}} p^{(i)}(s'^{(i)}|s^{(i)}, a^{(i)}) V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)}) \\
&= \sum_{s'^{(i)} \in \mathcal{S}^{(i)}} p^{(i)}(s'^{(i)}|s^{(i)}, a^{(i)}) \left[ r^{(i)}(s^{(i)}, s'^{(i)}) + \gamma V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)}) \right]
\end{aligned}
$$

Then we derive the value difference:

$$
V^{\pi^{(i)}}(s^{(i)}, \mathcal{T}^{(i)}) - V^{\pi^{(j)}}(s^{(i)}, \mathcal{T}^{(j)})
$$
$$
= \sum_{s'^{(i)}} p^{(i)}(s'^{(i)}|s^{(i)}, a^{(i)}) \left[ r^{(i)}(s^{(i)}, s'^{(i)}) + \gamma V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)}) \right] - \sum_{s'^{(j)}} p^{(j)}(s'^{(j)}|s^{(j)}, a^{(j)}) \left[ r^{(j)}(s^{(j)}, s'^{(j)}) + \gamma V^{\pi^{(j)}}(s'^{(j)}, \mathcal{T}^{(j)}) \right]
$$

<span style="color:blue">*minus and plus</span> $\sum_{s'^{(i)}} p^{(j)}(G(s'^{(i)})|s^{(j)}, a^{(j)}) \left[ r^{(i)}(s^{(i)}, s'^{(i)}) + \gamma V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)}) \right]$

$$
= \sum_{s'^{(i)}} p^{(i)}(s'^{(i)}|s^{(i)}, a^{(i)}) \left[ r^{(i)}(s^{(i)}, s'^{(i)}) + \gamma V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)}) \right] - \sum_{s'^{(i)}} p^{(j)}(G(s'^{(i)})|s^{(j)}, a^{(j)}) \left[ r^{(i)}(s^{(i)}, s'^{(i)}) + \gamma V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)}) \right]
$$
$$
+ \sum_{s'^{(i)}} p^{(j)}(G(s'^{(i)})|s^{(j)}, a^{(j)}) \left[ r^{(i)}(s^{(i)}, s'^{(i)}) + \gamma V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)}) \right] - \sum_{s'^{(j)}} p^{(j)}(s'^{(j)}|s^{(j)}, a^{(j)}) \left[ r^{(j)}(s^{(j)}, s'^{(j)}) + \gamma V^{\pi^{(j)}}(s'^{(j)}, \mathcal{T}^{(j)}) \right]
$$

<span style="color:blue">*combine first two terms, rewrite the third term given invertible function $G$</span>

$$
= \sum_{s'^{(i)}} \left[ p^{(i)}(s'^{(i)}|s^{(i)}, a^{(i)}) - p^{(j)}(G(s'^{(i)})|s^{(j)}, a^{(j)}) \right] \left[ r^{(i)}(s^{(i)}, s'^{(i)}) + \gamma V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)}) \right]
$$
$$
+ \sum_{s'^{(j)}} p^{(j)}(s'^{(j)}|s^{(j)}, a^{(j)}) \left[ r^{(i)}(G^{-1}(s^{(j)}), G^{-1}(s'^{(j)})) + \gamma V^{\pi^{(i)}}(G^{-1}(s'^{(j)}), \mathcal{T}^{(i)}) \right]
$$
$$
- \sum_{s'^{(j)}} p^{(j)}(s'^{(j)}|s^{(j)}, a^{(j)}) \left[ r^{(j)}(s^{(j)}, s'^{(j)}) + \gamma V^{\pi^{(j)}}(s'^{(j)}, \mathcal{T}^{(j)}) \right] \quad \text{<span style="color:blue">*combine last two terms, note the assumption of reward function</span>}
$$
$$
= \sum_{s'^{(i)}} \left[ p^{(i)}(s'^{(i)}|s^{(i)}, a^{(i)}) - p^{(j)}(G(s'^{(i)})|s^{(j)}, a^{(j)}) \right] \left[ r^{(i)}(s^{(i)}, s'^{(i)}) + \gamma V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)}) \right]
$$
$$
+ \gamma \sum_{s'^{(j)}} p^{(j)}(s'^{(j)}|s^{(j)}, a^{(j)}) \left[ V^{\pi^{(i)}}(G^{-1}(s'^{(j)}), \mathcal{T}^{(i)}) - V^{\pi^{(j)}}(s'^{(j)}, \mathcal{T}^{(j)}) \right]
$$

Therefore, the absolute value of value difference can be upper bounded. The proof is similar to the proof of Theorem 2.

$$
\begin{aligned}
\left| V^{\pi^{(i)}}(s^{(i)}, \mathcal{T}^{(i)}) - V^{\pi^{(j)}}(s^{(j)}, \mathcal{T}^{(j)}) \right| \quad &\leq \quad 2M D_{TV}(p^{(i)}(\cdot|s^{(i)}, a^{(i)}), p^{(j)}(G(\cdot)|s^{(j)}, a^{(j)})) \\
&+ \quad \gamma \sup_{s'^{(j)} \in \mathcal{S}_B^{(j)}} \left| V^{\pi^{(i)}}(G^{-1}(s'^{(j)}), \mathcal{T}^{(i)}) - V^{\pi^{(j)}}(s'^{(j)}, \mathcal{T}^{(j)}) \right| \\
&\leq \quad d + \gamma \sup_{s'^{(i)} \in \mathcal{S}_B^{(i)}} \left| V^{\pi^{(i)}}(s'^{(i)}, \mathcal{T}^{(i)}) - V^{\pi^{(j)}}(G(s'^{(i)}), \mathcal{T}^{(j)}) \right| \\
&\leq \quad \frac{d}{1 - \gamma}
\end{aligned}
$$

$\square$

Obviously, if the state equivalence relation is only true for identical states (i.e. $G$ is an identity mapping, $s^{(i)} B s^{(j)}$ if and only if $s^{(i)} = s^{(j)}$), then Proposition 2 degenerates into Proposition 1. If we optimize the action translator $H$ to minimize $d$ for policy $\pi^{(j)}$ and $\pi^{(i)}(s^{(i)}) = H(s^{(j)}, \pi^{(j)}(s^{(j)}))$, the policy value for correspondent states $s^{(i)}$ and $s^{(j)}$ can be close. Minimizing $d$ means finding actions leading to next states remaining correspondent.

### E.2 Method

According to Proposition 2, we not only learn an action translator $H$, but also state translators $G$ mapping target states $s^{(i)}$ to the equivalent states on source task $\mathcal{T}^{(j)}$ and $G^{-1}$ identifying correspondent state on target task $\mathcal{T}^{(i)}$. We additionally learn a discriminator network $D$ to assist learning of state translator.

Given transition data $s^{(j)}$ on source task and $s^{(i)}$ on target task, the adversarial objective is:

$$
\min_G \max_D \mathcal{L}_{adv}(G, D) = \log D(s^{(j)}) + \log(1 - D(G(s^{(i)})))
$$
.

$G$ aims to map target state $s^{(i)}$ to the distribution of states on source task, while $D$ tries to distinguish translated state $G(s^{(i)})$ and real states in the source task. To build state equivalence, the translated state should be translated back to the source state. We further leverage cycle consistency loss to learn the one-to-one mapping on states across tasks:

$$
\mathcal{L}_{back} = |G^{-1}(G(s^{(i)})) - s^{(i)}| + |G(G^{-1}(s^{(j)})) - s^{(j)}|
$$

Drawn upon Proposition 2, we extend our transfer loss $\mathcal{L}_{trans}$ to $\mathcal{L}_{trans,s,a}$. Formally,

$$
\mathcal{L}_{trans,s,a} = -\log F(\tilde{s}_{t+1}^{(i)} | \tilde{s}_t^{(i)}, \tilde{a}_t^{(i)})
$$

where $\tilde{s}_{t+1}^{(i)} = G^{-1}(s_{t+1}^{(j)})$, $\tilde{s}_t^{(i)} = G^{-1}(s_t^{(j)})$, and $\tilde{a}^{(i)} = H(s_t^{(j)}, a_t^{(j)})$. $\mathcal{L}_{trans,s,a}$ is applied to optimize the state translator $G^{-1}$ and action translator $H$.

In this way, given the state $s_t^{(j)}$ on source task, we first get the correspondent state $\tilde{s}_t^{(i)}$ on target task. Then the translated action $\tilde{a}^{(i)}$ make transition to next state $\tilde{s}_{t+1}^{(i)}$ on target task still correspondent to next state $s_{t+1}^{(i)}$ on source task. The objective function $\mathcal{L}_{trans,s,a}$ drives the next state distribution on the target task $p^{(i)}(\cdot|\tilde{s}_t^{(i)}, \tilde{a}^{(i)})$ to be close to the distribution of correspondent next state on the source task $p^{(j)}(G(\cdot)|s_t^{(j)}, a_t^{(j)})$. This is implicitly minimizing $d$ in Proposition 2.

In practice, we may need the action translator network $H$ or the state translator network $G$ and $G^{-1}$ reasonably initialized, in order to prevent the joint training collapsing to a trivial solution. The implementation details of learning the context model, forward dynamics model and action translator are the same as we explained in Appendix C.2. During training of the state translator, the weight of $\mathcal{L}_{adv}, \mathcal{L}_{back}, \mathcal{L}_{trans\_s\_a}$ is 10, 30, 100 respectively, the same as default hyper-parameters in (Zhang et al. 2020c). The similar technique of learning state translator and action translator has been mentioned in (Zhang et al. 2020c). Yet, our theorems shed light on its underlying mechanism and our objective function for learning the action translator is simpler.

### E.3 Experiments on Tasks Differing in Reward Function

When the tasks share **the same state space and action space but the reward function varies**, we combine our action translator with a state translator for policy transfer.

To investigate this scheme for policy transfer, we conduct experiments on MetaWorld task moving the robot arm to a goal location. We set the source and target task with different goal locations and hence with different reward functions. Tab. 15 lists the goal locations on the source and target tasks. Specifically, on the same state $s$ of the agent's current location $(x, y, z)$, the reward varies across tasks, because it is inversely proportional to the distance from the current location to goal. The initial location of robot arm is randomly sampled between $[-0.1, 0.6, 0.02]$ and $[0.1, 0.7, 0.02]$. The state is current location of the robot arm. The action is the moving vector of the robot arm.

We compare our method and (Zhang et al. 2020c) learning both state translator and action translator. We initialize the state translator networks by assigning $G(s = (x, y, z)) = G^{-1}(s = (x, y, z)) = (-x, y, z)$. As observed in Tab. 15, ours compares favorably with (Zhang et al. 2020c) and achieves satisfactory cumulative episode reward on the target task. We conclude that, for source and target tasks with different reward functions depending on the state and next state, learning state translator and action translator jointly is promising for policy transfer.

| Source Task | Target Task | Source policy on source task | Source policy on target task | Transferred policy (Zhang et al. 2020c) on target task | Transferred policy (Ours) on target task |
|---|---|---|---|---|---|
| $[-0.1, 0.8, 0.2]$ | $[0.1, 0.8, 0.2]$ | 4855.7 | 947.5 | $1798.2_{(\pm 592.4)}$ | $\mathbf{3124.3}_{(\pm 1042.0)}$ |
| $[-0.1, 0.8, 0.2]$ | $[0.05, 0.8, 0.2]$ | 4855.7 | 1470.2 | $1764.0_{(\pm 316.3)}$ | $\mathbf{1937.1}_{(\pm 424.5)}$ |
| $[-0.1, 0.8, 0.2]$ | $[0.1, 0.8, 0.05]$ | 4855.7 | 1040.8 | $\mathbf{2393.7}_{(\pm 869.8)}$ | $2315.7_{(\pm 1061.5)}$ |
| 2-leg | 3-leg | 5121.4 | NA | $1957.8_{(\pm 298.4)}$ | $\mathbf{2018.2}_{(\pm 50.8)}$ |

Table 15: Mean ($\pm$ standard error) of episode rewards over 3 runs, comparing source and transferred policy on target task. This is expanding Tab. 15 in the main text.

## E.4 Experiments on Tasks Differing in State and Action Space

For tasks with **different state space and action space**, we investigate the proposed idea on MuJoco environment HalfCheetah. The HalfCheetah agent by default has 2 legs in the source task and we modify the agent to have 3 legs in the target task. Because the agents have different numbers of joints in the source and target task, the dimensions of state space and action space also differ, as explained in (Zhang et al. 2020c). Again, we compare our method and (Zhang et al. 2020c) learning both state translator and action translator. We assign a good initialization for the action translator in both methods as (Zhang et al. 2020c) introduced. We remark that ours with a simpler objective function and fewer components than the baseline method can transfer the source policy to perform well on the target task.