

TRANSFER-CONTROLLABLE POLICY FOR MODEL PROTECTION IN DEEP REINFORCEMENT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Online deep reinforcement learning (DRL) suffers from sample inefficiency. This inefficiency challenges the training of effective policy models for complex tasks and demands substantial time and computing resources. As trained policy models can be transferred to other applications, protecting their intellectual property (IP) has become a pressing issue. To address this, we need to prevent unauthorized transfers for IP protection while maintaining transferability for future scalability. We propose the first Transfer-Controllable Reinforcement Learning (TCRL) framework. It has two key components: the Environment Randomization module generates unauthorized target-domain environments randomly, and the Transfer-Controllable module trains a policy model using source-domain and these unauthorized target-domain environments. This model resists transfer in unauthorized settings yet remains transferable in authorized ones. We validated the framework’s effectiveness across various DRL environments and algorithms. The TCRL policy model is hard to transfer to similar unauthorized target-domain environments, but achieves source-domain-like performance in authorized ones. In the MuJoCo environment, our trained policy model attains 98.78% of the source-domain performance in authorized target-domain environments, and only 50.38% in unauthorized ones.

1 INTRODUCTION

Deep Reinforcement Learning (DRL) techniques have thrived in various AI fields, like video games Nie et al. (2024), board games Schrittwieser et al. (2020), and robot control Han et al. (2024); Haarnoja et al. (2024). However, significant expertise is needed to ensure their proper operation Miki et al. (2022). For example, by creating 8 different reward functions, including torque and joint speed costs, and adopting curriculum learning, researchers enabled legged robots to learn animal-like dynamic maneuvers Hwangbo et al. (2019). Also, training the AlphaGo policy model requires tens of millions of dollars and thousands of GPUs Silver et al. (2016). Given the high investment in time, resources, and expertise, protecting the intellectual property (IP) of policy models is crucial.

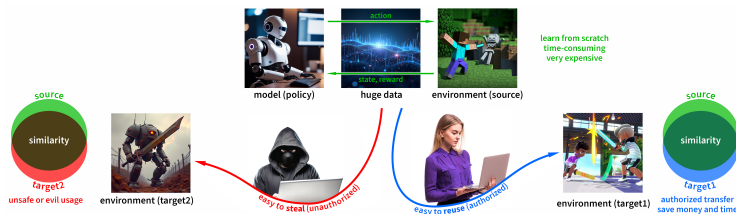


Figure 1: Training a policy model from scratch is time-consuming and costly. However, this model can be seamlessly transferred to similar scenarios, substantially reducing training expenses. To safeguard against theft by malicious actors, any transfer to unauthorized environments must be strictly prohibited. Simultaneously, to guarantee the model’s scalability in future applications, its transferability within authorized environments should be maintained.

DRL policy models risk theft and unauthorized transfer. Their relatively small model size (Fig. 1) facilitates easy theft and quick transfer to similar domains. During policy training Silver et al. (2016), these models learn from observations and generate actions, storing valuable knowledge, making

054 them more vulnerable to theft than large datasets. Training a policy model from scratch is extremely
 055 time-consuming and costly, while using a pre-trained model can boost efficiency. Since policy mod-
 056 els hold environment-related knowledge, it can be transferred to target domains via methods like
 057 learning from demonstrations, representation transfer, and inter-task mapping Yi et al. (2023). This
 058 transferability, however, also makes them prone to abuse.

059 Competitors may misuse obtained policy models by transferring them to similar scenarios, violating
 060 IP rights. For example, if a trained gameplay robot’s policy model leaks, it could be used for illegal
 061 activities like poaching through transfer learning, as depicted in Fig. 1. However, completely banning
 062 model transfer across different environments would harm the open-source community and limit
 063 legitimate applications. With the growing use of DRL techniques, protecting policy model IP has be-
 064 come an urgent issue. To address this, we propose the first **Transfer-Controllable Reinforcement**
 065 **Learning (TCRL)** framework. This framework aims to balance model IP protection and usability
 066 in authorized environments. It has two main modules: the Environment Randomization module,
 067 which randomly generates unauthorized target-domain environments, and the Transfer-Controllable
 068 Training module. The latter optimizes data from the source and authorized domains and performs
 069 reverse optimization on unauthorized target-environment data. We also design a new policy-model
 070 objective to stabilize the training process. In our experiment, the transfer difficulty of all environ-
 071 ments is set equally to ensure consistent experimental conditions. Our main contributions are as
 072 follows:

- 073 • We propose a new transfer-controllable task in DRL and validate its existence.
- 074 • We propose a preliminary TCRL framework to address this transfer-controllable task.
- 075 • Experimental results show policy models from our framework are controllably transferable:
 076 readily transferring to authorized target domains, yet struggling with unauthorized ones.
 077

078 2 RELATED WORKS

080 **Policy Transfer in DRL.** Our work is the opposite of the goal of policy transfer. Policy transfer
 081 uses the knowledge learned on the source domain to help the policy training on the target domain
 082 Zhu et al. (2023). In policy distillation, the algorithms learn a student policy π_{θ_S} by minimizing
 083 the divergence of action distributions between the teacher policy π_{θ_T} and the student policy π_{θ_S}
 084 according to trajectories τ . These studies can be further divided into two categories: teacher distil-
 085 lation Allen et al. (2021); Xu et al. (2019); Zhu et al. (2022) and student distillation D’Eramo et al.
 086 (2024); Schmitt et al. (2018). The difference between them is that τ is sampled from teacher policy:
 087 $\tau \sim \pi_{\theta_T}$ in teacher distillation and student policy: $\tau \sim \pi_{\theta_S}$ in student distillation. In policy reuse,
 088 the algorithms reuse a set of teacher policies by the means of π -reuse exploration strategy, which
 089 defines the trade-off among exploitation of the student policy, exploitation of the teacher policies,
 090 and exploration of random actions using the evaluation of the teacher policies’ performance on the
 091 target domain. The typical research include Wu et al. (2024); Daoudi et al. (2024); Zhang et al.
 092 (2024); Gimelfarb et al. (2021); Tao et al. (2021); Yi et al. (2023); Tian et al. (2023).

093 **IP Protection in Deep Learning.** The IP protection in DRL is still in its infancy, whereas research
 094 on IP protection in Supervised Learning (SL) has made significant progress. In SL, the research
 095 can be divided into three main categories: digital watermarking, backdoor and fingerprint Xue et al.
 096 (2021); Fkirin et al. (2022). Digital watermarking involves embedding robust digital watermarks into
 097 SL models to protect the model IP rights Uchida et al. (2017). The side effect of digital watermarking
 098 that reduces the model prediction abilities is optimized from two aspects by backdoor Adi et al.
 099 (2018) and fingerprint Zhao et al. (2020). In DRL, some attack techniques are proposed to change
 100 the model output Behzadan & Munir (2017); Chen et al. (2021b), which shows that it is urgent
 101 to study countermeasures of IP infringement on DRL models Ilahi et al. (2021). Similar to the
 102 SL methods, some research in DRL also embeds watermarks into the target policy for ownership
 verification Behzadan & Hsu (2019); Chen et al. (2021a).

103 Different from the watermarking-based methods above, transfer-controllable learning restricts the
 104 generalization ability of the model on target domains while preserving its performance on source
 105 domains. The first approach of non-transfer learning was proposed in SL Wang et al. (2022). How-
 106 ever, in DRL, the transfer-controllable learning problem has yet to be studied, and there are still
 107 many issues to be addressed in order to protect model IP. Compared to the large model size and
 stable training dataset in SL, the DRL model size is relatively small, and the dataset during training
 is unstable.

3 MOTIVATION

Unlike SL, policy model initialization in DRL is crucial Yi et al. (2023). Online DRL faces two major challenges: the exploration-exploitation dilemma and sparse rewards. The former requires balancing between using existing policies for rewards and exploring with stochastic policies; GoExplore addresses this by storing environmental states in an archive buffer Ecoffet et al. (2021). The latter occurs when agents need extended action sequences for non-zero rewards, which can be mitigated through immediate intrinsic rewards, as demonstrated with the 11 distinct rewards designed for a bipedal robot Duan et al. (2021). Addressing these challenges demands substantial resources in terms of funding, hardware, and training time.

However, if we have a better initial policy model before training, these difficulties can be alleviated Tirinzoni et al. (2019); Van Baar et al. (2019); Dennis et al. (2020); Abdolshah et al. (2021). A good initial model can perform correct actions, reducing the need for extensive exploration in the target environment to obtain sparse rewards Barreto et al. (2017); Wulfmeier et al. (2017); Riedmiller et al. (2018); Li et al. (2019); Guo et al. (2022).

Training a transfer-controllable policy model can resist transfer attacks and protect model IP rights. This raises two key questions: **(1) Is training such a model necessary?** **(2) What are the specific challenges in training transfer-controllable models in DRL compared to SL?** Given that a well-initialized policy model can reduce DRL training difficulty through transfer learning, answering these questions is significant.

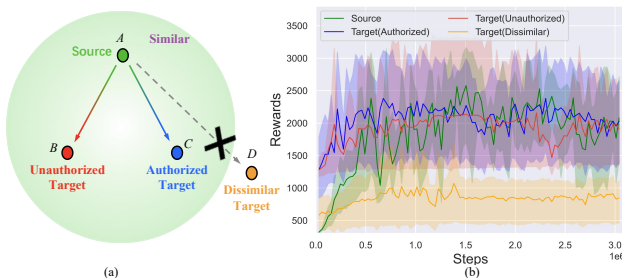


Figure 2: Preliminary Experiment. (a) We assume a green area exists where the source-domain policy model A can transfer to models B and C. Due to the policy overfitting in DRL, this green area is usually thought non-existent, meaning model A hardly transfers to model D. (b) The experimental results verify the existence of the green area. The Source, Target(Authorized), Target(Unauthorized), and Target(Dissimilar) curves correspond to policy models A, B, C, and D respectively.

Necessity of Training Transfer-controllable Model. To tackle Question (1), we first consider whether the policy model space contains similar regions. Such similarity is key as it enables a well-trained source-domain policy to transfer smoothly to certain target domains. In DRL, policies typically overfit to the source-domain environment, hindering their transfer to target domains. As shown in the dissimilar target domain in Fig. 2(a), transferring Policy A to Policy D is difficult. We assume there are similar target domains where the source-trained Policy A can quickly transfer to Policies B and C, as marked by the green area in Fig. 2(a).

To test our hypothesis, an experiment is conducted on the MuJoCo Hopper robot (Fig. 4(b)). The source domain featured Hopper parameters (torso, thigh, foot) of (0.05, 0.05, 0.06). Target domains, authorized, unauthorized, and dissimilar, have parameters (0.10, 0.05, 0.06), (0.05, 0.10, 0.06), and (0.2, 0.05, 0.06), respectively. Policy A, trained in the source domain, is transferred to these target domains. Results (Fig. 2(b)) show Policy A achieved 2000 reward in the source domain. In target domains with altered torso (Policy B) or thigh (Policy C) sizes, performance quickly reaches 2000. However, in the dissimilar target domain (Policy D) with a large torso change, Policy A’s overfitting to the source domain hinders transfer. Results show training a transfer-controllable policy is essential. The source-domain trained policy has some transferability. We must prevent its transfer to unauthorized domains, while ensuring transfer to authorized ones for future scalability.

Different Research Points on SL and DRL. For Question (2), the research interests of transfer-controllable learning technology in DRL are distinct from those in SL. In SL, the main problem is how to overfit the source domain model to limit its generalization ability on the target do-

main Sadashivaiah et al.. The parameter space of the SL model is large, thus providing many directions for its optimization, making it easier to control the direction of overfitting while still ensuring the model’s generalization on the target domain is limited. Furthermore, the datasets in SL are usually huge and stable, which makes the training process more stable and further reduces the difficulty of controlling the direction of overfitting. However, policy overfitting in DRL can limit the transfer of source domain policy models to certain target domain environments. However, in other target domain environments, the small model size of DRL models and the changing dataset distribution during its training process, bring more diverse problems in the DRL field. Therefore, it is necessary to conduct research on training transfer-controllable policy models.

4 METHODOLOGY

4.1 PRELIMINARY

In DRL, the agent learn from interaction with the environment, and the learning process is modeled with the Markov Decision Process (MDP) defined by a tuple (S, A, P, r, γ) . At each step t , the agent samples an action $a_t \in A$ from a policy distribution $\pi_\theta(a_t|s_t)$ where $s_t \in S$ is the observed state from the environment and θ is the policy model parameter. After passing the action a_t into the environment, the environment transmits into the next state s_{t+1} with the transition distribution $p(s_{t+1}|s_t, a_t) \in P$, and the agent receives a reward $r_t(s_t, a_t)$. Appendix A provides detailed explanations of each variable and foundational background on DRL.

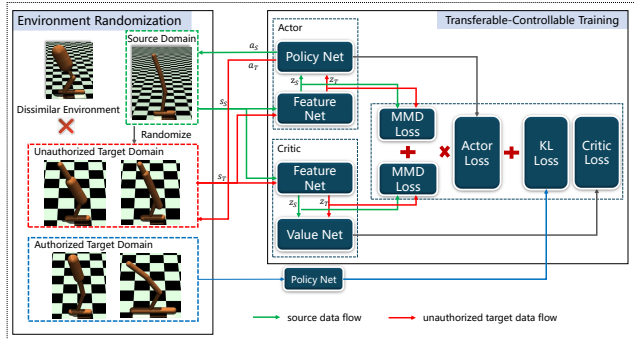


Figure 3: The main TCRL framework consists of Environment Randomization and Transfer-Controllable Training. The Environment Randomization module is used to randomly generate the unauthorized target domain environments and train some policy nets on authorized target domain environments, while the TCRL Training module trains the transfer-controllable policy model through a specific transfer-controllable loss function. The solid line represents the data interaction between the TCRL model and the environments, where the interaction targets of the red and green lines are the source domain and the unauthorized target domain environments, respectively.

4.2 TCRL FRAMEWORK

This paper introduces the TCRL framework (Fig. 3). The Environment Randomization module generates unauthorized target domains, uses user-provided authorized domains to train policies, and provides data for transfer-controllable training. The Transfer-Controllable Training module uses this output to train the transfer-controllable policy. Interactive source and unauthorized target data (Fig. 3) from source and generated target domains enable reverse transfer training, which limits transfer to unauthorized environments. Also, the authorized-domain policy uses KL divergence for scalability.

4.2.1 ENVIRONMENT RANDOMIZATION MODULE

As depicted in the left part of Fig. 3, the Environment Randomization module generates unauthorized target-domain environments and concurrently creates several authorized policy models based on the user-provided authorized target domain. Initially, it randomly selects source-domain policy models to fine-tune the authorized policy models and collects offline datasets during the fine-tuning process. Subsequently, it randomly generates unauthorized target-domain environments according to specified rules, such as the two robot environments within the red box. Next, through model fine-tuning, their transferability is evaluated. Dissimilar target environments, marked with a red cross

Algorithm 1 Environment Randomization Module

Input: environment parameters ρ and Actor model set $P_{\text{model}} = \{f_{\pi_{\theta}}g\}$ in source domain, parameter adjustment threshold δ , user-provided authorized target-domain environments E_{Auth}

Output: environment parameter set $E_{\text{Unauth.Target}} = \{f_{\rho_i}g_{i=0}^N\}$, authorized target-domain Actor models π^{Auth}

- 1: Randomly select π_{θ} from P_{model} ▷ Transfer authorized policy
- 2: Fine-tune π_{θ} on E_{Auth} to get π^{Auth}
- 3: Initialize $i = 0, E_{\text{Unauth.Target}} = \{f_{\rho}g\}$
- 4: **while** $i < N$ **do** ▷ Parameter randomization
- 5: Randomize parameters ρ_i in the range $[\rho - \delta, \rho + \delta]$
- 6: Construct E_i through parameters ρ_i
- 7: Randomly select π_{θ} from P_{model} ▷ Model fine-tuning
- 8: Fine-tune π_{θ} on E_i to get reward r_{target}
- 9: **if** Converge Time $t = T_{\text{threshold}}$ **then**
- 10: $E_{\text{Unauth.Target}} = E_{\text{Unauth.Target}} \cup \{f_{\rho_i}g, i = i + 1\}$
- 11: Calculate the scaling factor $f_r = r_{\text{target}}/r_{\text{source}}$
- 12: **end if** ▷ Screening unauthorized target environments
- 13: **end while**

in Fig. 3, are excluded because they deviate significantly from the source-domain environment. The objective of this paper is to obtain a source-domain policy model that is difficult to transfer in previously unauthorized target-domain environments. Therefore, this module randomly generates target-domain environments and selects those that are easily transferable. We derived Theorem 1 to elucidate the existence of such unauthorized environments in the target domain.

Theorem 1: Let τ_S and τ_T represent all optimal trajectories in the source and target domains, respectively. For a given δ , a state-action pair $(s_t, a_t, s_{t+1}) \in \tau_T$ is considered source-similar if there exists a state-action pair $(s_t^0, a_t^0, s_{t+1}^0) \in \tau_S$ such that $|j_{s_t} - s_t^0| < \delta$ and $|j_{s_{t+1}} - s_{t+1}^0| < \delta$. Conversely, a state-action pair is considered target-specific if it is not source-similar. Then, an increase in the number of target-specific state-action pairs makes it more difficult to transfer to the target domain environment, and the $H\Delta H$ distance between the source and target domains satisfies

$$d_{H\Delta H}(\tilde{D}_S, \tilde{D}_T) = 2 \sup_{\eta \in \mathcal{H}_d} |\Pr_{\mathcal{D}_S}[z : \eta(z) = 1] - \Pr_{\mathcal{D}_T}[z : \eta(z) = 1]| \quad (1)$$

where z denotes the feature of the state s , \tilde{D}_S and \tilde{D}_T represents the dataset on the source and target domain, respectively. The detailed proof for Theorem 1 is included in the Appendix B.

In detail, the algorithm process can be divided into four main phases: fine-tune authorized policy, parameter randomization, unauthorized model fine-tuning and screening environment, as shown in Algorithm 1. More details in the Appendix D.

4.2.2 TRANSFER-CONTROLLABLE TRAINING MODULE

The Transfer-Controllable Training module, illustrated in the right part of Fig. 3, is designed to train a transfer-controllable policy model in the source domain. This module interacts with the Environment Randomization module, as depicted in the middle of Fig. 3. During the model training process, the Actor model receives the source domain states s_S and the target domain states s_T from the environments in the Environment Randomization module at each step, and output the corresponding actions a_S and a_T . Subsequently, the specific policy model objective is defined as

$$J_{\text{TCRL}}^{\theta_k, D_{\text{Source}}, D_{\text{Unauth.Target}}, D_{\text{Auth.Target}}}(\theta) = J^{\theta_k, D_{\text{Source}}}(\theta) - \eta L_{\text{MMD}}^{\theta_k, D_{\text{Unauth.Target}}}(\theta) + \lambda (\hat{D}_{KL}^{D_{\text{Auth.Target}}}(\pi_{\theta}(j_{s_t})) \|\pi^{\text{Auth}}(j_{s_t}))) \quad (2)$$

where $J^{\theta_k, D}(\theta)$ is defined in Eq. (6) and Eq. (7), θ_k denotes the Actor model parameters after k th training, D_{Source} and $D_{\text{Unauth.Target}}$ are data buffers, η represents the learning rate of the reverse training, λ represents the weighting factors for authorized scalability, and \hat{D}_{KL} is the Kullback-Leibler divergence function. Meanwhile, the Feature Net in the Actor model outputs the intermediate features z_s and z_t , and the maximum mean discrepancy (MMD) loss is computed as

$$L_{\text{MMD}} = \min \left(\alpha, \beta \left\| \sum_{i=1}^{n_1} \Phi(z_{s,i}) - \sum_{i=1}^{n_2} \Phi(z_{t,i}) \right\|_H^2 \right) \quad (3)$$

Algorithm 2 Transfer-Controllable Training Module

Input: environment parameter set $E_{\text{Unauth.Target}} = \{f, g, D_{\text{Unauth.Target}}\}$ in target domain, maximum data buffer size jD
Output: transfer-controllable Actor model π_θ

- 1: Initialize $k = 0, D_{\text{Source}} = f, g, D_{\text{Unauth.Target}} = f, g$ ▷ Algorithm preparation
- 2: Randomize the parameters of Actor π_{θ_k} and Critic v_{ϕ_k}
- 3: Construct E_{Source} and $f, E_k, g_{k=1}^L$ through $E_{\text{Unauth.Target}}$ in each domain
- 4: **while** $k < N$ **do** ▷ Data collection
- 5: **while** $jD_{\text{Source}} + jD_{\text{Unauth.Target}} < jD$ **do**
- 6: Collect τ_S by running π_{θ_k} in source domain
- 7: Collect τ_T by running π_{θ_k} in authorized target domain
- 8: $D_{\text{Source}} = D_{\text{Source}} \cup \tau_S, g$
- 9: $D_{\text{Unauth.Target}} = D_{\text{Unauth.Target}} \cup \tau_T, g$
- 10: **end while** ▷ Auxiliary variable calculation
- 11: Compute $\hat{A}_t^{\text{Source}}$ and $\hat{A}_t^{\text{Unauth.Target}}$ with Eq. (5)
- 12: Compute \hat{R}_t on D_{Source} and $D_{\text{Unauth.Target}}$ with $\hat{R}_t = \hat{A}_t + v_{\phi_k}(s_t)$
- 13: **repeat** ▷ Model parameter update
- 14: Randomly choose (s_t, a_t, z_t) from datasets D
- 15: Recompute $\pi_\theta(a_t|s_t)$ and $v_\phi(s_t)$
- 16: Compute the MMD loss L_{MMD} with Eq. (3)
- 17: Update π_{θ_k} by maximizing $J_{\text{TCRL}}^{\theta_k}(\theta)$ through $\theta_{k+1} = \theta_k + \gamma_\theta J_{\text{TCRL}}^{\theta_k}(\theta)$
- 18: Update v_{ϕ_k} on $L_{\text{MSE}}(\phi)$ through $\phi_{k+1} = \phi_k + \gamma_\phi L_{\text{MSE}}(\phi)$
- 19: **until** D_{Source} is empty
- 20: $k = k + 1, D_{\text{Source}} = f, g, D_{\text{Unauth.Target}} = f, g$
- 21: **end while**

where $\Phi(\cdot)$ denotes the Gaussian kernel function, H indicates the Hilbert space, and α, β are the tunable hyperparameters.

The two equations above are essential for achieving anti-transfer training in the unauthorized target-domain environments while maintaining transferability in the authorized target-domain environments. In Eq. (2), the first term represents the model training in the source domain environment, while the second term indicates the reverse model training in the generated unauthorized target domain environments. There may be some similar samples on the source domain dataset D_{Source} and the target domain dataset $D_{\text{Unauth.Target}}$. This causes the gradients from the source domain $J^{\theta_k, D_{\text{Source}}}(\theta)$ and the target domain $J^{\theta_k, D_{\text{Auth.Target}}}(\theta)$ to be opposite, negatively impacting the model training on the source domain environment. To address this, factors η and L_{MMD} are introduced to adjust the strength of reverse training on the target domain environments, thus decreasing the negative impact. Additionally, the term L_{MMD} could increase the distribution distance between the source domain feature z_s and the target domain feature z_t , making it easier to optimize in different directions on D_{Source} and $D_{\text{Unauth.Target}}$, thus reducing the difficulty of reverse training optimization. In addition, the third term $\hat{D}_{KL}^{D_{\text{Auth.Target}}}$ is used to ensure the transferability of the policy model in the authorized target-domain environment. Here, $D_{\text{Auth.Target}}$ is the fixed dataset obtained in the previous step, which is used to fine-tune the transfer-controllable policy model to the policy model in the authorized target domain. Since $D_{\text{Auth.Target}}$ is a fixed and small dataset, it has little impact on the overall training of the first two terms. From the domain adaptation theory, we derived Theorem 2 to illustrate the role of the MMD loss as follows

Theorem 2: Assume $p(s, a)$ is the joint distribution of state s and action a . Given $\delta \in [0, 1]$, let a partition $\Omega \subset \mathbb{R}^n$ on the H space satisfies $P_{p(s,a)}(s \in \Omega) = \delta$, then

(1) there exists a partition Ω_{D_S} and Ω_{D_T} such that

$$d_{H(H)}(D_S, D_T) \geq 2|E_{s \sim D_S}[A(s) \notin A^\theta(s)] - E_{s \sim D_T}[A(s) \notin A^\theta(s)]| \quad (4)$$

(2) maximizing the MMD loss is equivalent to increasing the distance $d_{H(H)}$. The detailed proof for Theorem 2 is included in the Appendix B.

Concretely, the data processing flow includes four main phases: preparation, data collection, auxiliary variable calculation, and model parameter update, as shown in Algorithm 2. More details in the

Appendix D. The discounted reward \hat{R}_t is calculated as

$$\hat{A}_t^D = \sum_{D,l} (\gamma\lambda)^l (r_t + \gamma v_{\phi_k}(s_{t+l+1}) + \beta L_{\text{MMD}}(\log(\pi(a_t/s_t)) \mathbf{1}_a + \epsilon) v_{\phi_k}(s_{t+l})) \quad (5)$$

where D denotes the data buffer, v_{ϕ_k} denotes the Critic model, γ and λ are adjustment factors.

5 EXPERIMENTAL RESULT

5.1 EXPERIMENT SETUP

To verify the training effect of our framework and the performance of TCRL model obtained through training, we conducted experiments on different DRL algorithms and in different test environments.

DRL Algorithms, namely DQN Mnih et al. (2015) and PPO Schulman et al. (2017), are employed to comprehensively evaluate the performance of these algorithms under various conditions, aiming to uncover their respective advantages and limitations in solving the targeted problems.

Test Environments. The main body of the text primarily presents the experimental results of our framework in the Maze Environment and the MuJoCo Environment Todorov et al. (2012). The configuration examples of these environments are illustrated in Fig. 4. In the experiments conducted in the Maze Environment, we test the effectiveness of the DQN algorithm within our framework. Meanwhile, in the experiments carry out in the MuJoCo Environment, we examine the performance of the PPO algorithm within the same framework. Additional experimental results under various settings are available in Appendix E.

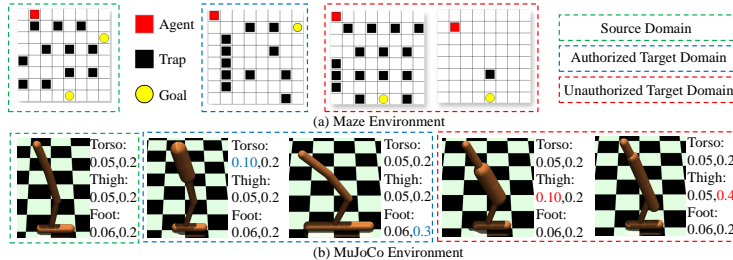


Figure 4: Overview of Experiment Setup. (a) Maze Environment. It consists of the Agent, Trap, and Goal. In the source domain, there are two Goals, one on the right and one at the bottom. The authorized target domain has a single Goal on the right, while the unauthorized target domain has only one Goal at the bottom. In independent experiments, the positions of the Agent, Trap, and Goal vary; (b) MuJoCo Environment. It encompasses MuJoCo robots with diverse configurations. In the authorized target domain, users set the configurations based on the subsequent scalability requirements of the model. In contrast, configurations in the unauthorized target domain are randomly generated by the Environment Randomization module. In this example, users primarily specify the Torso and Foot configurations of the Hopper robot, while the Thigh configuration is generated by the Environment Randomization module.

5.2 PERFORMANCE OF DQN ON MAZE ENVIRONMENT

In this experiment, Agent receives a final reward of 60 upon reaching Goal and -10 if it enters Trap by mistake. A single experiment terminates when Agent reaches the Goal or the environment runs for more than 200 steps. As shown in Fig. 5(a), both the original DQN (blue curve) and TCRL_DQN (orange curve) can achieve a reward value of around 50 during training in the source domain, indicating that the trained Agents can complete the tasks. This implies that our method has little impact on the performance of the source-domain policy model during training.

Fig. 5(b) reveals that in the policy model transfer experiment, the transfer-controllable policy model trained by the TCRL framework (orange curve) can complete the task in the authorized environment but struggles to do so in the unauthorized environment. Here, TCRL_Trans_Unauth_1 (green curve) and TCRL_Trans_Unauth_2 (red curve) correspond to the two experimental settings in the red box of Fig. 4(a) respectively. The red curve shows that even when the number of Traps is significantly reduced, the transfer-controllable policy model still fails to complete the task. This demonstrates

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431



Figure 5: Experiment of DQN on Maze Environment.

that the policy model obtained through our training exhibits strong reverse transfer ability in the unauthorized environment.

5.3 TRAINING PERFORMANCE OF PPO ON MUJoCo ENVIRONMENT

In this experiment, we aim to verify the effect of the TCRL algorithm on the training performance of the baseline algorithm. To do so, we used 32 copies of the same source domain environment to train the benchmark PPO algorithm in parallel, and employ the same 32 source domain environment copies, as well as 32 unauthorized target domain environments, to train the TCRL algorithm in parallel. The convergence of the reward curve is used as the evaluation criteria.

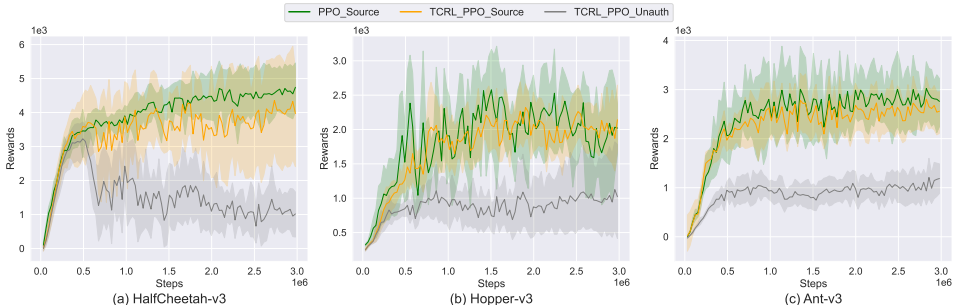


Figure 6: Training performance of the baseline PPO algorithm and our TCRL algorithm. The blue *PPO_Source* and orange *TCRL_PPO_Source* curves denote the performance of PPO and TCRL on the source domain, while the green *TCRL_PPO_Unauth* curves indicate the performance of TCRL on the unauthorized target domains.

The experimental results show that during training, TCRL can achieve a performance similar to PPO in the source domain. Meanwhile, it can significantly degrade the performance of the policy model in the unauthorized target domain. In Fig. 6, for the HalfCheetah, Hopper, and Ant tasks, the TCRL_PPO_Source curve converges to an average reward value close to that of the PPO.Source curve, though with a slightly larger variance. This indicates that our method may slightly increase the training difficulty of the algorithm, but has minimal impact on the final training outcome, as both can yield effective policy models. On the other hand, the TCRL_PPO_Unauth curve is limited to a very low value in the unauthorized target domain environment. This demonstrates that our method restricts the policy model’s performance in such environments, laying the groundwork for subsequent transfer experiments.

5.4 TRANSFERRING PERFORMANCE OF PPO ON MUJoCo ENVIRONMENT

In this experiment, we aim to verify the effectiveness of the obtained transfer-controllable policy model in preventing the transfer of the source domain to the unauthorized target domain. To do so, we use the trained PPO and TCRL policy models to transfer on 8 authorized and 32 unauthorized target domain environments, respectively. Additionally, a random initialized policy model was trained under the same target domain environment as a benchmark. Subsequently, the policy models were tested on 8 authorized and 8 unauthorized target domain environments to verify the average transfer-controllable ability during training process. The convergence of the reward curve was then used as the evaluation criterion for transfer performance.

Based on the experimental results, it is evident that the TCRL policy model can effectively impede the transferability of the source-domain policy model to the unauthorized target domain. As depicted in Fig. 7, the reward values achieved by the TCRL_PPO_Trans_Unauth curve are substantially lower



Figure 7: Comparing the transfer performance between the PPO and TCRL policy models. The *PPO_Trans_Auth* curve denotes the PPO model and the *TCRL_PPO_Trans_Auth* curve denotes the TCRL model on authorized target domain, while the *PPO_Trans_Unauth* curve denotes the PPO model and the *TCRL_PPO_Trans_Unauth* curve denotes the TCRL model on unauthorized target domain. The green *PPO_Random* curve, trained with a random initialized model, serves as the baseline.

Table 1: Transferring performance on MuJoCo environment. *PPO_Random* denotes the average rewards in authorized and unauthorized target-domain environments. *TCRL_Trans_Auth* and *TCRL_Trans_Unauth* represent the average rewards of TCRL policy model transfer to authorized and unauthorized target-domain environments, respectively.

	HalfCheetah-v3	Hopper-v3	Ant-v3	Mean
PPO_Random	4123.66	2057.22	2598.95	-
TCRL_Trans_Auth	4207.33	2075.46	2427.83	-
Ratio	102.03%	100.89%	93.42%	98.78%
TCRL_Trans_Unauth	2516.65	1028.16	1043.02	-
Ratio	61.03%	49.98%	40.13%	50.38%

than those of the *PPO_Trans_Unauth* curve. This implies that the TCRL policy model encounters significant difficulties in migrating to the unauthorized target-domain environment, thereby demonstrating a robust anti-transfer capacity. The *PPO_Random* curve represents the average reward values obtained through training from the initial state in each task environment. The convergence values of the *PPO_Trans_Unauth* curve are comparable to those of the *PPO_Random* curve. This indicates that the original PPO algorithm is essentially incapable of preventing the source-domain policy model from transferring to the unauthorized target-domain environment.

Meanwhile, Fig. 7 clearly demonstrates that TCRL policy model preserves its transferability within the authorized target domain, thereby providing an avenue for subsequent model expansion. The convergence values of the *TCRL_PPO_Trans_Auth* curve exhibit minimal divergence from those of the *PPO_Trans_Auth* curve and closely approximate the reward values of the *PPO_Random*. This observation implies that our proposed policy model retains a high-level of transferability.

As indicated in Table 1, the TCRL policy model derived from our training regimen not only sustains a transfer performance of 98.78% in the authorized target domain but also effectively restricts the transfer performance of the policy model to the unauthorized target domain to 50.38%.

6 CONCLUSION AND LIMITATION

In this paper, we have introduced a new task of training transfer-controllable policies in DRL and presented an original framework to address this task. Firstly, we have examined the necessity of transfer-controllable learning in DRL and identified potential challenges that may arise. Subsequently, we proposed the TCRL framework for transfer-controllable training and theoretically demonstrated its feasibility. Moreover, we applied this framework to obtain a transfer-controllable policy model and empirically validated its efficacy in safeguarding against transfer attacks on the policy model. However, the TCRL framework’s major limitation is that it consumes approximately twice the computational resources of conventional DRL training, mainly because of the high computational cost of stochastically generating suitable unauthorized target domain.

REFERENCES

- 486
487
488 Majid Abdolshah, Hung Le, Thommen Karimpanal George, Sunil Gupta, Santu Rana, and Svetha
489 Venkatesh. A new representation of successor features for transfer across dissimilar environments.
490 In *International Conference on Machine Learning*, pp. 1–9. PMLR, 2021.
- 491 Yossi Adi, Carsten Baum, Moustapha Cisse, Benny Pinkas, and Joseph Keshet. Turning your weak-
492 ness into a strength: Watermarking deep neural networks by backdooring. In *27th USENIX Security
493 Symposium (USENIX Security 18)*, pp. 1615–1631, 2018.
- 494
495 Cameron Allen, Neev Parikh, Omer Gottesman, and George Konidaris. Learning markov state ab-
496 stractions for deep reinforcement learning. *Advances in Neural Information Processing Systems*,
497 34:8229–8241, 2021.
- 498
499 André Barreto, Will Dabney, Rémi Munos, Jonathan J Hunt, Tom Schaul, Hado P van Hasselt,
500 and David Silver. Successor features for transfer in reinforcement learning. *Advances in neural
501 information processing systems*, 30, 2017.
- 502
503 Vahid Behzadan and William Hsu. Sequential triggers for watermarking of deep reinforcement
504 learning policies. *arXiv preprint arXiv:1906.01126*, 2019.
- 505
506 Vahid Behzadan and Arslan Munir. Vulnerability of deep reinforcement learning to policy induction
507 attacks. In *International Conference on Machine Learning and Data Mining in Pattern Recognition*,
508 pp. 262–275. Springer, 2017.
- 509
510 Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wort-
511 man Vaughan. A theory of learning from different domains. *Machine learning*, 79:151–175,
512 2010.
- 513
514 Kangjie Chen, Shangwei Guo, Tianwei Zhang, Shuxin Li, and Yang Liu. Temporal watermarks
515 for deep reinforcement learning models. In *Proceedings of the 20th International Conference on
516 Autonomous Agents and MultiAgent Systems*, pp. 314–322, 2021a.
- 517
518 Kangjie Chen, Shangwei Guo, Tianwei Zhang, Xiaofei Xie, and Yang Liu. Stealing deep reinforce-
519 ment learning models for fun and profit. In *Proceedings of the 2021 ACM Asia Conference on
520 Computer and Communications Security*, pp. 307–319, 2021b.
- 521
522 Paul Daoudi, Bogdan Robu, Christophe Prieur, Ludovic Dos Santos, and Merwan Barlier. Enhancing
523 reinforcement learning agents with local guides. *arXiv preprint arXiv:2402.13930*, 2024.
- 524
525 Michael Dennis, Natasha Jaques, Eugene Vinitzky, Alexandre Bayen, Stuart Russell, Andrew Critch,
526 and Sergey Levine. Emergent complexity and zero-shot transfer via unsupervised environment
527 design. *Advances in neural information processing systems*, 33:13049–13061, 2020.
- 528
529 Carlo D’Eramo, Davide Tateo, Andrea Bonarini, Marcello Restelli, and Jan Peters. Sharing knowl-
530 edge in multi-task deep reinforcement learning. *arXiv preprint arXiv:2401.09561*, 2024.
- 531
532 Helei Duan, Jeremy Dao, Kevin Green, Taylor Apgar, Alan Fern, and Jonathan Hurst. Learning task
533 space actions for bipedal locomotion. In *2021 IEEE International Conference on Robotics and
534 Automation (ICRA)*, pp. 1276–1282. IEEE, 2021.
- 535
536 Adrien Ecoffet, Joost Huizinga, Joel Lehman, Kenneth O Stanley, and Jeff Clune. First return, then
537 explore. *Nature*, 590(7847):580–586, 2021.
- 538
539 Alaa Fkirin, Gamal Attiya, Ayman El-Sayed, and Marwa A Shouman. Copyright protection of deep
540 neural network models using digital watermarking: a comparative study. *Multimedia Tools and
541 Applications*, 81(11):15961–15975, 2022.
- 542
543 Michael Gimelfarb, Scott Sanner, and Chi-Guhn Lee. Contextual policy transfer in reinforcement
544 learning domains via deep mixtures-of-experts. In *Uncertainty in Artificial Intelligence*, pp. 1787–
545 1797. PMLR, 2021.
- 546
547 Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.

- 540 Yijie Guo, Qiucheng Wu, and Honglak Lee. Learning action translator for meta reinforcement
541 learning on sparse-reward tasks. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
542 pp. 6792–6800, 2022.
- 543
- 544 Tuomas Haarnoja, Ben Moran, Guy Lever, Sandy H Huang, Dhruva Tirumala, Jan Humplik, Markus
545 Wulfmeier, Saran Tunyasuvunakool, Noah Y Siegel, Roland Hafner, et al. Learning agile soccer
546 skills for a bipedal robot with deep reinforcement learning. *Science Robotics*, 9(89):eadi8022,
547 2024.
- 548 Lei Han, Qingxu Zhu, Jiapeng Sheng, Chong Zhang, Tingguang Li, Yizheng Zhang, He Zhang,
549 Yuzhen Liu, Cheng Zhou, Rui Zhao, et al. Lifelike agility and play in quadrupedal robots using
550 reinforcement learning and generative pre-trained models. *Nature Machine Intelligence*, 6(7):
551 787–798, 2024.
- 552
- 553 Jemin Hwangbo, Joonho Lee, Alexey Dosovitskiy, Dario Bellicoso, Vassilios Tsounis, Vladlen
554 Koltun, and Marco Hutter. Learning agile and dynamic motor skills for legged robots. *Science
555 Robotics*, 4(26):eaau5872, 2019.
- 556 Inaam Ilahi, Muhammad Usama, Junaid Qadir, Muhammad Umar Janjua, Ala Al-Fuqaha, Dinh Thai
557 Hoang, and Dusit Niyato. Challenges and countermeasures for adversarial attacks on deep rein-
558 forcement learning. *IEEE Transactions on Artificial Intelligence*, 3(2):90–109, 2021.
- 559
- 560 Siyuan Li, Rui Wang, Minxue Tang, and Chongjie Zhang. Hierarchical reinforcement learning
561 with advantage-based auxiliary rewards. *Advances in Neural Information Processing Systems*,
562 32, 2019.
- 563
- 564 Xingyu Liu, Deepak Pathak, and Kris M Kitani. Revolver: Continuous evolutionary models for
565 robot-to-robot policy transfer. In *International Conference on Machine Learning*, 2022.
- 566 Takahiro Miki, Joonho Lee, Jemin Hwangbo, Lorenz Wellhausen, Vladlen Koltun, and Marco Hut-
567 ter. Learning robust perceptive locomotion for quadrupedal robots in the wild. *Science Robotics*,
568 7(62):eabk2822, 2022.
- 569
- 570 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Belle-
571 mare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level
572 control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- 573
- 574 Buqing Nie, Jingtian Ji, Yangqing Fu, and Yue Gao. Improve robustness of reinforcement learn-
575 ing against observation perturbations via lipschitz policy networks. In *Proceedings of the AAAI
576 Conference on Artificial Intelligence*, volume 38, pp. 14457–14465, 2024.
- 577
- 578 Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel
579 Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement
580 learning and demonstrations. *arXiv preprint arXiv:1709.10087*, 2017a.
- 581
- 582 Aravind Rajeswaran, Kendall Lowrey, Emanuel V. Todorov, and Sham M Kakade. Towards gen-
583 eralization and simplicity in continuous control. In *Advances in Neural Information Processing
584 Systems*, volume 30, 2017b.
- 585
- 586 Martin Riedmiller, Roland Hafner, Thomas Lampe, Michael Neunert, Jonas Degraeve, Tom Wiele,
587 Vlad Mnih, Nicolas Heess, and Jost Tobias Springenberg. Learning by playing solving sparse
588 reward tasks from scratch. In *International conference on machine learning*, pp. 4344–4353.
589 PMLR, 2018.
- 590
- 591 Vijay Sadashivaiah, Keerthiram Murugesan, Ronny Luss, Pin-Yu Chen, Chris Sims, James Hendler,
592 and Amit Dhurandhar. To transfer or not to transfer: Suppressing concepts from source represen-
593 tations. *Transactions on Machine Learning Research*.
- 594
- 595 Simon Schmitt, Jonathan J Hudson, Augustin Zidek, Simon Osindero, Carl Doersch, Wojciech M
596 Czarnecki, Joel Z Leibo, Heinrich Kuttler, Andrew Zisserman, Karen Simonyan, et al. Kickstart-
597 ing deep reinforcement learning. *arXiv preprint arXiv:1803.03835*, 2018.

- 594 Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon
595 Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering atari,
596 go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.
- 597
- 598 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
599 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 600
- 601 David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche,
602 Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering
603 the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- 604 Yunzhe Tao, Sahika Genc, Jonathan Chung, Tao Sun, and Sunil Mallya. Repaint: Knowledge
605 transfer in deep reinforcement learning. In *International Conference on Machine Learning*, pp.
606 10141–10152. PMLR, 2021.
- 607
- 608 Zikang Tian, Ruizhi Chen, Xing Hu, Ling Li, Rui Zhang, Fan Wu, Shaohui Peng, Jiaming Guo,
609 Zidong Du, Qi Guo, et al. Decompose a task into generalizable subtasks in multi-agent rein-
610 forcement learning. In *Thirty-seventh Conference on Neural Information Processing Systems*,
611 2023.
- 612
- 613 Andrea Tirinzoni, Mattia Salvini, and Marcello Restelli. Transfer of samples in policy search via
614 multiple importance sampling. In *International Conference on Machine Learning*, pp. 6264–
6274. PMLR, 2019.
- 615
- 616 Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control.
617 In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 5026–5033.
618 IEEE, 2012. doi: 10.1109/IROS.2012.6386109.
- 619
- 620 Yusuke Uchida, Yuki Nagai, Shigeyuki Sakazawa, and Shin’ichi Satoh. Embedding watermarks
621 into deep neural networks. In *Proceedings of the 2017 ACM on international conference on
622 multimedia retrieval*, pp. 269–277, 2017.
- 623
- 624 Jeroen Van Baar, Alan Sullivan, Radu Cordorel, Devesh Jha, Diego Romeres, and Daniel Nikovski.
625 Sim-to-real transfer learning using robustified controllers in robotic tasks involving complex dy-
626 namics. In *2019 International Conference on Robotics and Automation (ICRA)*, pp. 6001–6007.
627 IEEE, 2019.
- 628
- 629 Lixu Wang, Shichao Xu, Ruiqi Xu, Xiao Wang, and Qi Zhu. Non-transferable learning: A new
630 approach for model ownership verification and applicability authorization. In *International Con-
631 ference on Learning Representations*, 2022.
- 632
- 633 Chengjie Wu, Pingzhong Tang, Jun Yang, Yujing Hu, Tangjie Lv, Changjie Fan, and Chongjie
634 Zhang. Conservative offline policy adaptation in multi-agent games. *Advances in Neural Infor-
635 mation Processing Systems*, 36, 2024.
- 636
- 637 Markus Wulfmeier, Ingmar Posner, and Pieter Abbeel. Mutual alignment transfer learning. In
638 *Conference on Robot Learning*, pp. 281–290. PMLR, 2017.
- 639
- 640 Keyulu Xu, Jingling Li, Mozhi Zhang, Simon S Du, Ken-ichi Kawarabayashi, and Stefanie Jegelka.
641 What can neural networks reason about? *arXiv preprint arXiv:1905.13211*, 2019.
- 642
- 643 Mingfu Xue, Yushu Zhang, Jian Wang, and Weiqiang Liu. Intellectual property protection for deep
644 learning models: Taxonomy, methods, attacks, and evaluations. *IEEE Transactions on Artificial
645 Intelligence*, 1(01):1–1, 2021.
- 646
- 647 Qi Yi, Rui Zhang, Shaohui Peng, Jiaming Guo, Yunkai Gao, Kaizhao Yuan, Ruizhi Chen, Siming
Lan, Xing Hu, Zidong Du, et al. Online prototype alignment for few-shot policy transfer. *arXiv
preprint arXiv:2306.07307*, 2023.
- Gengzhi Zhang, Liang Feng, Yu Wang, Min Li, Hong Xie, and Kay Chen Tan. Reinforcement
learning with adaptive policy gradient transfer across heterogeneous problems. *IEEE Transactions
on Emerging Topics in Computational Intelligence*, 2024.

648 Yuchen Zhang, Tianle Liu, Mingsheng Long, and Michael Jordan. Bridging theory and algorithm
649 for domain adaptation. In *International conference on machine learning*, pp. 7404–7413. PMLR,
650 2019.

651 Jingjing Zhao, Qingyue Hu, Gaoyang Liu, Xiaoqiang Ma, Fei Chen, and Mohammad Mehedi Has-
652 san. Afa: Adversarial fingerprinting authentication for deep neural networks. *Computer Commu-
653 nications*, 150:488–497, 2020.

654 Jinhua Zhu, Yingce Xia, Lijun Wu, Jiajun Deng, Wengang Zhou, Tao Qin, Tie-Yan Liu, and
655 Houqiang Li. Masked contrastive representation learning for reinforcement learning. *IEEE Trans-
656 actions on Pattern Analysis and Machine Intelligence*, 45(3):3421–3433, 2022.

657 Zhuangdi Zhu, Kaixiang Lin, Anil K Jain, and Jiayu Zhou. Transfer learning in deep reinforcement
658 learning: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.

661 A USE OF LLMs

662 This paper did not use any LLMs during the research, writing, and other related processes.

663 B FOUNDATIONAL BACKGROUND ON DRL

664 B.1 DQN AND PPO OBJECTIVES AND CRITIC LOSS

665 In this paper, both the deep Q -network (DQN) and the proximal policy optimization (PPO) Schul-
666 man et al. (2017) algorithms are used to train the policy model. The main objectives of DQN and
667 PPO are:

$$668 J_{\text{DQN}}^{\theta_k, D_k}(\theta) = \hat{E}_{t, D_k} \left\{ \left[Q(s_t, a_t; \theta) - r_t(s_t, a_t) - \gamma \max_{a^0} Q(s_{t+1}, a^0; \theta_k) \right]^2 \right\}, \quad (6)$$

669 and

$$670 J_{\text{PPO}}^{\theta_k, D_k}(\theta) = \hat{E}_{t, D_k} \left\{ \min \left[\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_k}(a_t | s_t)} \hat{A}_t, \text{clip} \left(\frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_k}(a_t | s_t)}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_t \right] \right\}, \quad (7)$$

671 where Q represents the action-value function, θ_k indicates the network parameters of the old pol-
672 icy model at the k th training epoch, D_k denotes the data buffer at the k th training epoch, ϵ is a
673 hyperparameter, and \hat{A}_t indicates the advantage estimates.

674 The Critic loss L_{MSE} is defined as

$$675 L_{\text{MSE}}(\phi) = \hat{E}_{t, D_k} [(V_{\phi}(s_t) - \hat{R}_t)^2] \quad (8)$$

676 where D_k indicates the data buffer of the chosen (s_t, a_t, z_t) pairs from D_{Source} and $D_{\text{Unauth-Target}}$.

677 B.2 SYMBOL DEFINITIONS

678 The symbols used in this paper and their corresponding meanings are shown in Table below.

679 C THEORY PROOFS

680 **Theorem 1:** Let τ_{D_S} and τ_{D_T} represent all optimal trajectories in the source and target domains,
681 respectively. For a given δ , a state-action pair $(s_t, a_t, s_{t+1}) \in \tau_{D_T}$ is considered source-similar if
682 there exists a state-action pair $(s_t^0, a_t^0, s_{t+1}^0) \in \tau_{D_S}$ such that $|j s_t - s_t^0| < \delta$ and $|j s_{t+1} - s_{t+1}^0| < \delta$.
683 Conversely, a state-action pair is considered target-specific if it is not source-similar. Then, an
684 increase in the number of target-specific state-action pairs makes it more difficult to transfer to the
685 target domain environment, and the $H\Delta H$ distance between the source and target domains satisfies

$$686 d_{H\Delta H}(\tilde{D}_S, \tilde{D}_T) \leq 2 \sup_{\eta} j \Pr_{D_S}[z : \eta(z) = 1] - \Pr_{D_T}[z : \eta(z) = 1] \quad (9)$$

Table 2: Symbol Definitions

Symbol	Notation
t	The current time step
k	The k th training epoch
$S \subseteq \mathbb{R}^m$	The state space
$A \subseteq \mathbb{R}^n$	The action space
H	The Hilbert space
$P : S \times A \rightarrow \mathcal{P}(S)$	The state transition distribution
$r : S \times A \rightarrow \mathbb{R}$	The reward function
$\gamma \in [0, 1]$	The discounted factor
$s_t \in S$	The observed state from the environment at time step t
$a_t \in A$	The agent action at time step t
\mathcal{P}	The partition on the H space
z_t, f_t	The feature of the state s_t
$p(s_{t+1} s_t, a_t) \in \mathcal{P}$	The transition distribution at time step t
$r_t(s_t, a_t)$	The environment reward at time step t
(s_t, a_t, s_{t+1})	The state-action pair at time step t
θ	The policy model parameter
ϕ	The value network parameter
θ_k	The policy model parameter at k th training epoch
ϕ_k	The value network parameter at k th training epoch
$\pi_\theta(a_t s_t)$	The policy distribution at time step t with model parameter θ
$r_t(\theta)$	The policy probability ratio with model parameter θ
\mathcal{D}_k	The data buffer at the k th training epoch
\hat{A}_t	The advantage estimates at time step t
$J^{\theta_k, \mathcal{D}_k}(\theta)$	The main optimized objective of the PPO and the DQN algorithm
$J_{\text{TCRL}}^{\theta_k, \mathcal{D}_{\text{Source}}, \mathcal{D}_{\text{Unauth.Target}}}(\theta)$	The specific policy model optimized objective of the TCRL algorithm
L_{MMD}	The maximum mean discrepancy loss
τ_S, τ_T	The optimal trajectories in the source and target domains, respectively
$\mathcal{D}_S, \mathcal{D}_T$	The dataset on the source and target domain, respectively
$\mathcal{D}_{\text{Source}}$	The data buffer on the source domain
$\mathcal{D}_{\text{Unauth.Target}}$	The data buffer on the unauthorize target domain
$\mathcal{D}_{S, T}$	The partition on source and target domain dataset, respectively
ρ	The environment parameters
ϵ, δ	The hyperparameters representing small values
α, β	The tunable hyperparameters
$\min(\cdot)$	The minimize function
$\text{clip}(\cdot)$	The clip function
$E(\cdot)$	The expected function
$U(\cdot)$	The uniform distribution
$\mathcal{K}(\cdot)$	The Gaussian kernel function
$\text{Pr}(\cdot)$	The probability function

where z denotes the feature of the state s , \tilde{D}_S and \tilde{D}_T represents the dataset on the source and target domain, respectively.

Proof: Firstly, the RL transfer problem needs to be transformed into an SL optimization problem.

Assume that a trajectory τ is randomly selected from the set of target domain trajectories τ_{D_T} . If any state-action pair $(s_t, a_t, s_{t+1}) \in \tau$ is source-similar, the optimal actions a_t and a_t^0 satisfies that $|j_{a_t} - j_{a_t^0}| < \Delta$ as the state assumption conditions that $|j_{s_t} - j_{s_t^0}| < \delta$ and $|j_{s_{t+1}} - j_{s_{t+1}^0}| < \delta$. That is, if the state-action pairs on the target domain are all source-similar, then these optimal actions a_t and a_t^0 can be divided into different categories.

Moreover, if a state-action pair $(s_t, a_t, s_{t+1}) \in \tau$ is target-specific, suppose that $|j_{s_t} - j_{s_t^0}| < \delta$ and $|j_{s_{t+1}} - j_{s_{t+1}^0}| < 2\delta$, then there exists a state-action pair $(s_t^0, a_t^0, s_{t+1}^0)$ satisfies that $|j_{s_t^0} - j_{s_t^0}| < \delta$ and $|j_{s_{t+1}^0} - j_{s_{t+1}^0}| < 2\delta$. Then, the optimal actions satisfies that $|j_{a_t^0} - j_{a_t^0}| < \Delta$ and $|j_{a_t^0} - j_{a_t^0}| < \Delta$, and it means that $|j_{a_t} - j_{a_t^0}| < 2\Delta$. Furthermore, if there are few target-specific points in the target domain, these optimal actions a_t and a_t^0 can be divided into different categories through the auxiliary action a_t^{00} .

For the optimal trajectories $\tau_{D_S}(s_i) = [s_0, a_0^{opt}, \dots, a_i^{opt}, s_i]$, given the Markov property, optimizing τ_{D_S} in the source domain is equivalent to the existence of a classifier from state s_i to action a_i as $A_{D_S}(s_i) = a_i^{opt}$. Similarly, for the target domain, the optimal trajectories τ_{D_T} is equivalent to the optimal classifier $A_{D_T}(s_i) = a_i^{opt}$ from state s_i to action a_i .

Then, we derive the $H\Delta H$ distance between the source and target domain.

Let the action space be A . Since the action category space A_{D_S} and A_{D_T} are subsets of the action space A , and both the source domain classifier A_{D_S} and the target domain classifier A_{D_T} satisfy

$$A_{D_S} \cap A_{D_S} = A \quad \text{and} \quad A_{D_T} \cap A_{D_T} = A \quad (10)$$

the attribute of A_{D_S} and A_{D_T} is the same. Considering the network architecture of the policy model π_θ , assume the feature extraction function f_{D_S} of the Feature Net satisfies $f_{D_S}(s_i) = z_i \in \mathbb{R}^m$. As all classification problems can be transformed into binary classification Goodfellow et al. (2016), only the binary categories will be taken into consideration as $h : \mathbb{Z} \rightarrow \{0, 1\}$. Based on the domain adaptation theory Ben-David et al. (2010), for the classifier $A = h \circ f$, the error of the given classifier $h(z)$ on the target domain D_T is

$$\epsilon_{D_T}(h) < \epsilon_{D_S}(h) + \frac{1}{2} d_{H \cap H}(\tilde{D}_S, \tilde{D}_T) + \lambda \quad (11)$$

where $\epsilon_{D_T}(h)$ and $\epsilon_{D_S}(h)$ denote the error of the given classifier $h(z)$ on the source and target domain, respectively. The variable $d_{H \cap H}$ represents the generalized distance between data buffer \tilde{D}_S and \tilde{D}_T on the specific H space. Meanwhile, the const parameter λ satisfies that

$$\lambda = \epsilon_{D_S}(h^*) + \epsilon_{D_T}(h^*), \quad h^* = \arg \min_{h \in \mathcal{H}} \epsilon_{D_S}(h) + \epsilon_{D_T}(h) \quad (12)$$

where h^* indicates the best classifier with the lowest error sum λ of the source error ϵ_{D_S} and the target error ϵ_{D_T} on the H space. Meanwhile, the space $H\Delta H$ satisfies

$$H\Delta H = \{h : \eta(z^*) = 1\} \quad (13)$$

where define the variable z^* as

$$z^* = \{z : h_1(z) \oplus h_2(z), h_1, h_2 \in \mathcal{H}\} \quad (14)$$

where \oplus indicates the XOR operator.

Therefore, regarding the problem of transferring the source domain policy model π_θ into the target domain, it is equivalent to minimizing variables $\epsilon_{D_S}(h)$ and $d_{H \cap H}$. For minimizing the generalized distance $d_{H \cap H}$, we derive as follows

$$\begin{aligned} d_{H \cap H}(\tilde{D}_S, \tilde{D}_T) &= 2 \sup_{h_1, h_2 \in \mathcal{H}} |\Pr_{\tilde{D}_S}[fz : h_1(z) \oplus h_2(z)g] - \Pr_{\tilde{D}_T}[fz : h_1(z) \oplus h_2(z)g]| \\ &= 2 \sup_{\eta \in \mathcal{H}} |\Pr_{\tilde{D}_S}[fz : \eta(z) = 1g] - \Pr_{\tilde{D}_T}[fz : \eta(z) = 1g]| \\ &= 2 \sup_{\eta \in \mathcal{H}_d} |\Pr_{\tilde{D}_S}[z : \eta(z) = 1] - \Pr_{\tilde{D}_T}[z : \eta(z) = 1]| \end{aligned} \quad (15)$$

where H_d denotes the trained classifier space such that $h_1, h_2 \in H$.

Besides, as the number of target-specific state-action pairs increases, the difficulty of transferring the policy model from the source domain to the target domain increases from a geometric multiple. According to the Generalization Bound theorem Zhang et al. (2019), we have

$$\begin{aligned} \epsilon_{D_T}(f) &< \epsilon_{D_S}^{(\rho)}(f) + d_{f,F}^{(\rho)}(\tilde{D}_S, \tilde{D}_T) + \lambda + 2\sqrt{\frac{\log \frac{2}{\delta}}{2n}} + \sqrt{\frac{\log \frac{2}{\delta}}{2m}} \\ &+ \frac{2k^2}{\rho} \mathfrak{R}_{n,D_S}(\Pi_1 F) + \frac{2k}{\rho} \mathfrak{R}_{n,D_S}(\Pi_H F) + \frac{2k}{\rho} \mathfrak{R}_{m,D_T}(\Pi_H F) \end{aligned} \quad (16)$$

where f denotes all scoring functions, k represents the number of categories for classification problems in the source and target domains, ρ is a given const parameter, λ is a constant independent of f , \mathfrak{R} represents the Rademacher complexity, and $\Pi_1 F$ is defined as

$$\Pi_1 F = \{f(x, y) \mid y \in Y, f \in F\} \quad (17)$$

It can be seen from the above theorem that the increase of the number of categories k will lead to the increase of generalization error in the target domain. In our derivation, more target-specific state-action pairs mean more classification of action categories in both the source and target domains. That is to say, as the difference between the source domain and the target domain becomes larger, the generalization error between the source domain and the target domain will continue to increase. Furthermore, we can get that the increase of target-specific state-action pairs will make it more difficult for the policy model to transfer from the source domain to the target domain environment.

Theorem 2: Assume $p(s, a)$ is the joint distribution of state s and action a . Given $\delta \in [0, 1]$, let a partition $\Omega \subseteq \mathbb{R}^n$ on the H space satisfies $P_{p(s,a)}(s \in \Omega) = \delta$, then

(1) there exists a partition Ω_{D_S} and Ω_{D_T} such that

$$d_{H-H}(D_S, D_T) \leq 2 \left| E_{s \sim D_S}[A(s) \notin A^0(s)] - E_{s \sim D_T}[A(s) \notin A^0(s)] \right| \quad (18)$$

(2) maximizing the MMD loss is equivalent to increasing the distance d_{H-H} .

Proof: First, we prove that there exists a large upper bound of d_{H-H} that satisfies the transfer learning constraints for the transfer error ϵ_{D_T} .

Let $p(s, a)$ be the joint distribution of the state s and action a . A partition $\Omega \subseteq \mathbb{R}^n$ is constructed such that all states s in this partition Ω satisfies that

$$P_{p(s,a)}(s \in \Omega) = \delta \quad (19)$$

where the variable $\delta \in [0, 1]$. Given a classifier h , a classification method $k(z) = 1$ is generated on it, where $z \in \mathbb{R}^n$ and $h(z) > 0.5$. When $\delta = 1$, the partition Ω uniquely corresponds to a classifier k . In this case, the generalization error is

$$\epsilon(\Omega) = \epsilon(k) = E_j[A(s) \neq k(s)] \quad (20)$$

The optimal partition of the probability distribution p is denoted as

$$\Omega_p^* = \arg \min_{\Omega \subseteq \mathbb{R}^n} \epsilon(\Omega) \quad (21)$$

For the transfer problem on the target domain D_T , it is equivalent to the optimization problem

$$\min_{f,h} \epsilon_{D_S}(h \circ f), \quad \text{s.t. } f(s_{D_S}) = f(s_{D_T}) \quad (22)$$

For the transferred classifier A^{tran} , it belongs to the set of classifiers A^* that satisfy

$$\epsilon_{D_S}(h_{A^{\text{tran}}} \circ f_{A^{\text{tran}}}) \leq \epsilon_{D_S}(\Omega_{D_S}^*) \quad (23)$$

$$f_{A^{\text{tran}}}(s_{D_S}) = f_{A^{\text{tran}}}(s_{D_T}) \quad (24)$$

Consider the feature function $f(s)$ defined as follows: for a given partition Ω ,

$$f(s) = \begin{cases} 1_m, & s \in (S_{D_S} \setminus \Omega_{D_S}^*) \cup (S_{D_T} \setminus \Omega) \\ 0_m, & \text{otherwise} \end{cases} \quad (25)$$

where the parameter m denotes the dimensions of the feature vector. Let the classifier be $h(1_m) = 1$. Obviously, $A = h \circ f \circ s \in \mathcal{A}^*$. Construct that

$$\hat{\Omega} = \operatorname{argmax}_{\mathbb{R}^n} \epsilon_{D_T}(\Omega) \quad \text{s.t.} \quad P_{D_T}(s \in \Omega) = P_{D_S}(s \in \Omega_{D_S}^*) \quad (26)$$

the generalization error of the classifier

$$\hat{A} = h \circ f \circ s \in \mathcal{A}^* \quad (27)$$

corresponding to this partition is

$$\epsilon_{D_T}(\hat{A}) = \max_{A \in \mathcal{A}^*} \epsilon_{D_T}(A) \quad (28)$$

Define that

$$S^{\text{same}} = \{s \in S_{D_S} \setminus S_{D_T}\} \quad \text{and} \quad S^{\text{di}} = \{s \in S_{D_S} \setminus S_{D_T}\} \quad (29)$$

Assume that

$$P_{D_T}(s \in \Omega_{D_T}^*) = P_{D_S}(s \in \Omega_{D_S}^*) = 0.5 \quad (30)$$

and $s \in S^{\text{di}} \setminus \hat{\Omega}$, this approach still achieves optimization of source domain error while mapping D_S and D_T to the same distribution. In this case, it holds that

$$\max_{A \in \mathcal{A}^*} \epsilon_{D_T}(A) = (1 - \frac{|S^{\text{same}}|}{|S_{D_S} \cup S_{D_T}|}) (1 - \epsilon_{D_T}(\Omega_{D_T}^*)) + \epsilon_{D_T}(\Omega_{D_T}^*) \quad (31)$$

When $S^{\text{same}} = \emptyset$, it degenerates to

$$\max_{A \in \mathcal{A}^*} \epsilon_{D_T}(A) = 1 - \epsilon_{D_T}(\Omega_{D_T}^*) \quad (32)$$

This implies the existence of worst-case solutions that satisfy the original transfer learning conditions.

A worst-case classifier can be constructed as follows: Let Ω_{D_S} and Ω_{D_T} be chosen such that s has an equal probability of occurring in both the source and target domains and $S^{\text{same}} = \emptyset$. Define the feature function

$$f(s) = 1_m \quad \text{if } s \in (S_{D_S} \setminus \Omega_{D_S}) \cup (S_{D_T} \setminus \Omega_{D_T}) \quad (33)$$

$$f^\circ(s) = 1_m \quad \text{if } s \in (S_{D_S} \setminus \Omega_{D_S}) \cup (S_{D_T} \setminus (\mathbb{R}^n \cap \Omega_{D_T})) \quad (34)$$

and let the classifier be $h(1_m) = 1$. For the classifiers $A = h \circ f$ and $A^\circ = h \circ f^\circ$, both belong to classifiers that satisfy the transfer conditions, but there exists a $H\Delta H$ lower bound of

$$d_{H\Delta H} = 2|E_{D_S}[A(s) \neq A^\circ(s)] - E_{D_T}[A(s) \neq A^\circ(s)]| = 2 \quad (35)$$

the maximum value of $H\Delta H$ is achieved in this case.

Next, we aim to prove that increasing the MMD leads to an increase in the transfer error ϵ_{D_T} . The MMD distance is defined as

$$\begin{aligned} \text{MMD}(X, Y) &= \sqrt{\frac{1}{n} \sum_i \phi(x_i) \cdot \frac{1}{m} \sum_j \phi(y_j)} k_H \\ &= \sqrt{\frac{1}{n^2} \sum_i \sum_{i^\circ} \phi(x_i) \phi(x_{i^\circ}) + \frac{2}{nm} \sum_i \sum_j \phi(x_i) \phi(y_j) + \frac{1}{m^2} \sum_j \sum_{j^\circ} \phi(y_j) \phi(y_{j^\circ})} k_H \\ &= \sqrt{\frac{1}{n^2} \sum_i \sum_{i^\circ} k(x_i, x_{i^\circ}) + \frac{2}{nm} \sum_i \sum_j k(x_i, y_j) + \frac{1}{m^2} \sum_j \sum_{j^\circ} k(y_j, y_{j^\circ})} k \\ &= \sqrt{E(k(x_i, x_{i^\circ})) + 2E(k(x_i, y_j)) + E(k(y_j, y_{j^\circ}))} k \end{aligned} \quad (36)$$

where $x_i \in X$ and $y_j \in Y$, and the Gaussian kernel function is

$$k(u, v) = e^{-\frac{\|u-v\|^2}{\sigma^2}} \quad (37)$$

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

Consider these extreme scenarios:

i) When fixing $E(k(x_i, x_i^o)) = 1$ and $E(k(y_j, y_j^o)) = 1$, maximizing the MMD is equivalent to setting $E(k(x_i, y_j)) = 0$. By using the kernel function $k(u, v)$, we have

$$E(k(x_i, y_j)) = e^{-\frac{E(kx_i - y_j k^2)}{\sigma}} = 0 \quad (38)$$

which is equivalent to $E(kx_i - y_j k^2) \rightarrow +\infty$. Furthermore,

$$\begin{aligned} E(kx_i - y_j k^2) &= E(kx_i k^2 - 2kx_i k y_j k + ky_j k^2) \\ &= E(kx_i k^2) - 2E(kx_i k y_j k) + E(ky_j k^2) \\ &= E^2(kx_i k) - 2E(kx_i k)E(ky_j k) - 2k \text{Cov}(X, Y)k + E^2(ky_j k) \\ &= (E(kx_i k) - E(ky_j k))^2 - 2k \text{Cov}(X, Y)k \\ &\quad + 1 \end{aligned} \quad (39)$$

This is equivalent to that $k\bar{x} - \bar{y}k \rightarrow +\infty$.

ii) When fixing $E(k(x_i, y_j)) = 0$, maximizing the MMD is equivalent to setting $E(k(x_i, x_i^o)) = 1$ and $E(k(y_j, y_j^o)) = 1$. As before, this is equivalent to $E(kx_i - x_i^o k) \rightarrow 0$ and $E(ky_j - y_j^o k) \rightarrow 0$. Without loss of generality, we can assume that $kx_i k = kx_i^o k$ for $x_i, x_i^o \in X$. We consider the following on X :

$$\begin{aligned} E(kx_i - x_i^o k^2) &= E(kx_i k^2 - 2kx_i k x_i^o k + kx_i^o k^2) \\ &= E_i(E_{i^o}(kx_i k^2) - 2E_{i^o}(kx_i k x_i^o k) + E_{i^o}(kx_i^o k^2)) \\ &= E_i(kx_i k^2 - 2kx_i k E_{i^o}(kx_i^o k) + E_{i^o}(kx_i^o k^2)) \\ &= E(kx_i k^2 - 2\bar{x}kx_i k + E(kx_i k^2)) \\ &= 2(E(X^2) - E^2(X)) \\ &= 2D(X) \\ &\rightarrow 0 \end{aligned} \quad (40)$$

This is equivalent to $D(X) \rightarrow 0$ and $D(Y) \rightarrow 0$.

In summary, when optimizing the MMD, as it approaches the limit, we have

$$\lim_{\text{MMD} \rightarrow \max} \bar{x} - \bar{y} = +\infty \quad (41)$$

$$\lim_{\text{MMD} \rightarrow \max} D(x) = 0 \quad (42)$$

$$\lim_{\text{MMD} \rightarrow \max} D(y) = 0 \quad (43)$$

Considering the properties of limits, it is necessary that there exists a real number λ such that when $\text{MMD} > \lambda$, $\bar{x} - \bar{y}$ increases monotonically and $D(x)$ and $D(y)$ decrease monotonically. This means that there is a critical step after which the MMD training always descends the gradient towards the optimization of $\bar{x} - \bar{y}$, $D(x)$, and $D(y)$.

Considering with the Equation (31), when fixing other conditions and only considering the increase of $\bar{f}(s_{D_S}) - \bar{f}(s_{D_T})$, it is equivalent to a decrease in jS^{same}_j , which leads to an increase in $\epsilon_{D_T}(A)$.

Consider the feature extraction function $f(s) = 1_m$ for a given partition, where $s \in (S_{D_S} \setminus \Omega_{D_S}^* - (S_{D_T} \setminus \Omega))$. When fixing other conditions and considering the decrease of $D(f(s_{D_S}))$ and $D(f(s_{D_T}))$, we consider the conditions $(S_{D_S} \setminus \Omega_{D_S}^* - (S_{D_T} \setminus \Omega_{D_T}^*))$ and $(S_{D_S} \setminus \Omega_{D_S}^* - (S_{D_T} \setminus (R^n \cap \Omega_{D_T}^*)))$. To minimize the variance and achieve the optimal partition in the source domain, while ensuring that $j\Omega_{D_T}^* \setminus \Omega_{D_S}^*j$ approaches $jR^n \cap \Omega_{D_T}^* \setminus R^n \cap \Omega_{D_S}^*j$, the positive samples in the source domain and negative samples in the target domain are constrained to a point in the feature space. Similarly, this is also true for the negative samples in the source domain and positive samples in the target domain. Therefore, there exist only the optimal classifiers for D_S and D_T respectively in this feature space, and there does not exist a classifier that is optimal for both domains. Moreover, the partition boundary between the source and target classifiers is orthogonal.

In a word, maximizing the MMD loss is equivalent to increasing the distance $d_{H \perp H}$.

D IMPLEMENTATION DETAILS

D.1 NETWORK ARCHITECTURE

To build the Actor and Critic models, we use a three-layer MLP structure on the MuJoCo environment. The first two MLP layers act as feature extractors, while the last MLP layer is used as either the Policy Net or Value Net. The first two MLP layers are followed by a tanh activation function layer. The output of the last MLP layer of the Actor model is the mean value of the output policy, and the output of the last MLP layer of the Critic model is the estimated value of the current state.

D.2 HYPER PARAMETERS

In the Environment Randomization module, the scale parameter c is set to 1.5 for body mass, body inertia, and geom friction, and 1.3 for dof damping in the MuJoCo environment. For the tunable hyperparameters ϵ_1 and ϵ_2 are set to 0.1, ϵ_3 , and ϵ_4 , ϵ_1 , ϵ_2 , and ϵ_3 are set to 0.5 for each experiment, while ϵ_4 is set to 1 for the HalfCheetah-v3 and Hopper-v3 experiments and 3 for the Ant-v3 experiment. The tunable hyperparameter τ is set to 0.7 for the HalfCheetah-v3 and Ant-v3 experiments, and 0.8 for the Hopper-v3 experiment.

In the Transfer-Controllable Training module, the learning rate of the normal training is set to $3e-4$, and the learning rate of the reverse training is set to $3e-5$ for each experiment. The total buffer size is set to 4096, with the source domain dataset and the target domain dataset each being 2048, respectively. The step per epoch is set to 30000, and the step per collect is set to 2048. The batch size is set to 64, and the repeat per collect is set to 10. The thread number for collecting data is set to 64 during the model training process.

For the PPO algorithm, we employ both reward normalization and observation normalization techniques. In the loss function, the value function coefficient is set to 0.25, the entropy coefficient is set to 0.0, and the GAE lambda parameter is set to 0.95. Additionally, the epsilon clip parameter is set to 0.2.

D.3 UNAUTHORIZED TARGET DOMAIN ENVIRONMENTS ON TRAINING PROCESS

In the experiment, we use 32 threads to collect the state-action pair data on the source domain environment, while using 32 threads to collect the corresponding data on the unauthorized target domain environment. In order to ensure the diversity of data collected on the target domain environments, every 4 threads collect the state-action pair data obtained on the target domain environments with the same parameter configuration in parallel. In the subsequent supplementary experimental results section, we demonstrate that using this configuration for can achieve better training performance.

E ALGORITHM DETAILS

E.1 ENVIRONMENT RANDOMIZATION

As indicated in the main text of the paper, the process of Algorithm 1 can be divided into four main phases: fine-tune authorized model, parameter randomization, unauthorized model fine-tuning and screening environment.

In the phase of fine-tuning the authorized model, randomly select one from the source-domain policy models and transfer it to the authorized target domain environment given by the user, as shown in lines 1-2 of Algorithm 1.

In the parameter randomization phase, an unauthorized target domain environment is generated according to some custom randomization rules, as shown in lines 5-6 in Algorithm 1. This phase is mainly to randomize the relevant parameters in the source domain environment according to the characteristics of the source domain environment, so as to obtain the target domain environments similar to the MDP of the source domain environment.

In the unauthorized model fine-tuning phase, an Actor model π_θ is randomly selected from the source domain model set P_{model} for transfer learning, as shown in lines 7-8 in Algorithm 1. Then,

we retrain the Actor model initialized by π_θ in the generated target domain environments. In the MuJoCo environment used in the verification of this paper, this simple fine-tune method has been able to verify the availability of our framework. In practice, the appropriate transfer learning algorithm could be selected according to the environment characteristics.

In the screening environment phase, the unauthorized target domain environments are selected through the given custom rules, as shown in lines 10-11 in Algorithm 1. In this paper, we use the converge time and the source domain rewards to judge whether a target domain environment is easy to transfer. It is a simple and effective method to select the suitable target domain environments.

E.2 TRANSFER-CONTROLLABLE TRAINING

As indicated in the main text of the paper, the data processing flow of Algorithm 2 includes four main phases: algorithm preparation, data collection, auxiliary variable calculation, and model parameter update.

In the preparation phase, it mainly completes the construction of the environments and the initialization of related variables, as shown in lines 1-3 in Algorithm 2. The initial parameters of the target domain environments come from the result of Algorithm 1.

In the data collection phase, the Actor model π_{θ_k} is used to collect trajectories on the source domain environment E and the target domain environments $\{E_k\}_{k=1}^L$ respectively, as shown in lines 5-10 in Algorithm 2. When the sum of the capacities of the data buffer D_{Source} and $D_{\text{Unauth.Target}}$ is greater than the maximum threshold jD_j , the data collection phase ends.

In the auxiliary variable calculation phase, the discounted reward \hat{R}_t and the advantage estimates \hat{A}_t required for the subsequent training phase are calculated, as shown in lines 11-12 in Algorithm 2. Calculate \hat{A}_t^D using the generalized advantage estimation method on the data buffers D_{Source} and $D_{\text{Unauth.Target}}$.

In the model parameter update phase, the four loss functions, shown in Fig. 3, are used to update the model parameters, and the specific process is shown in lines 15-23 in Algorithm 2. In lines 14-15, the preparations before model training is completed. Then the MMD loss L_{MMD} and the Actor loss $J_{\text{TCRL}}^{\theta_k}(\theta)$ are calculated through Eq. (4) and Eq. (3), respectively. Next, the Critic loss is calculated through Eq. (9). Finally, the model parameters of the Actor network π_{θ_k} and the Critic network v_{ϕ_k} are updated using the gradient ascent method, as shown in lines 16-18.

F SUPPLEMENTARY EXPERIMENTAL RESULTS

F.1 ABLATION STUDIES ON EACH COMPONENT

To comprehensively evaluate the contribution of each key component within our proposed TCRL framework, we conducted additional ablation studies on the HalfCheetah-v3 benchmark. The results, presented in Table 3, quantify the impact of environment randomization and the transfer-controllable training module (specifically, the MMD loss).

Table 3: Ablation Studies on Environment Randomization and Transfer-Controllable Training Module. **w/o Env Filtering** refers to the variant where the process of screening and excluding unauthorized target environments during training is removed. **w/o MMD** indicates the removal of the Maximum Mean Discrepancy (MMD) loss from the policy model’s objective function.

Reward	w/o Env Filtering	w/o MMD	TCRL (full)
Unauthorized	1918	3012	2516
Authorized	3098	3985	4207

The experimental results highlight the significance of both components:

Impact of Environment Filtering: When Environment Filtering is omitted (“w/o Env Filtering”), the model’s performance degrades substantially. The reward in authorized scenarios drops from 4207 (TCRL full) to 3098. Concurrently, the reward in unauthorized scenarios is the lowest at

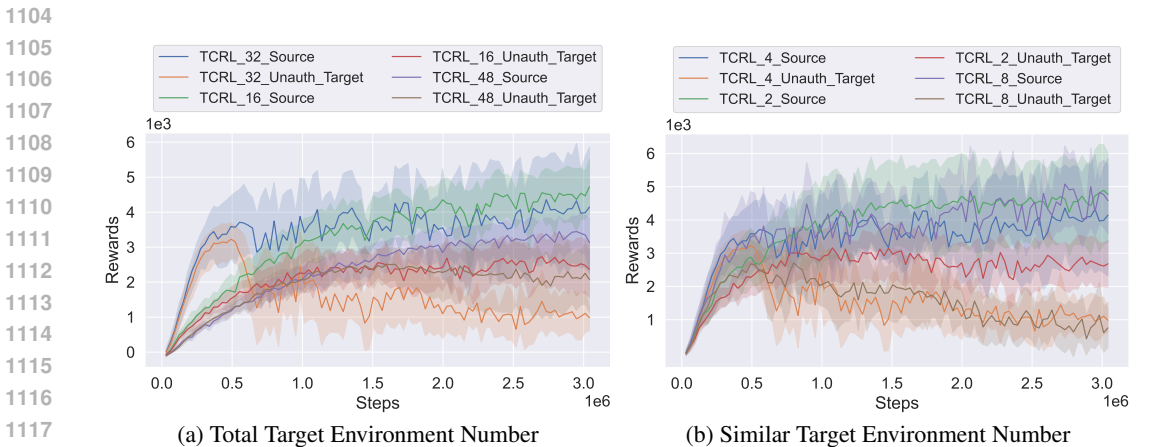
1080 1918. This outcome aligns with our hypothesis: without environment filtering, the training process
 1081 is exposed to unauthorized domains that may include dissimilar target environments. Such expo-
 1082 sure negatively impacts the model’s ability to learn an effective policy for authorized tasks and to
 1083 generalize appropriately.

1084 **Impact of MMD Loss:** Removing the MMD loss (“w/o MMD”) while retaining environment filter-
 1085 ing also leads to a noticeable performance decline compared to the full TCRL model. The authorized
 1086 reward decreases to 3985 from 4207, and the unauthorized reward increases to 3012 from 2516.
 1087 The MMD loss is designed to encourage the policy to learn domain-invariant representations of
 1088 state-action pairs, thereby helping to distinguish and adapt behaviors between authorized and unau-
 1089 thorized domains. Without it, the model struggles to effectively capture these crucial state-action
 1090 differences, leading to suboptimal performance in authorized settings and increased undesirable be-
 1091 havior in unauthorized ones.

1092 In contrast, the TCRL (full) model, which integrates both Environment Filtering and the MMD loss,
 1093 achieves the highest reward (4207) in authorized environments while maintaining a comparatively
 1094 lower reward (2516) in unauthorized environments. This demonstrates the synergistic effect of these
 1095 components in enabling robust and controllable transfer learning.

1097 F.2 DIFFERENT TARGET DOMAIN ENVIRONMENT CONFIGURATIONS ON TCRL TRAINING

1098
 1099 In this experiment, we mainly verify the impact of different unauthorized target domain environ-
 1100 ment number configurations on model performance during the training process. It mainly includes
 1101 changes in the total number of authorized target domain environments and changes in the proportion
 1102 of environments with the same configuration in the total target domain environments. In all these
 1103 experiments, we use 32 threads to collect data from the source domain.



1117
 1118 Figure 8: Training performance of different unauthorized target domain environment config-
 1119 urations on the HalfCheetah-v3 environment. The variables $TCRL_x_Source$ and
 1120 $TCRL_x_Unauth_Target$ represent the training performance of the TCRL algorithm on the source
 1121 and unauthorized target domains, respectively. (a) In this experiment, x refers to the number of
 1122 threads utilized for data collection in the target domain environments. Specifically, every 4 threads
 1123 were assigned to use the unauthorized target domain environments having an identical configuration.
 1124 (b) In this experiment, x denotes that every group of x threads was allocated to collect data from a
 1125 target domain environment having the identical configuration. In all these experiments, we use 32
 1126 threads to collect data on the source domain.

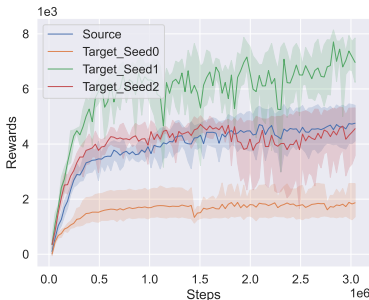
1127
 1128 In the Total Target Environment Number experiment, we change the total number of target do-
 1129 main environments, as shown in the left part of Fig. 8. The $TCRL_{32_Source}$ curve and the
 1130 $TCRL_{32_Unauth_Target}$ curve represent the default TCRL algorithm training configuration. Compar-
 1131 ing the $TCRL_{32_Unauth_Target}$ curve with the $TCRL_{16_Unauth_Target}$ curve, it can be seen
 1132 that reducing the total number of target domain environments will increase the reward value achieved
 1133 by TCRL in the target domain, which also means that reducing the total number of target domain en-
 1134 vironments reduces the effectiveness of TCRL algorithm in suppressing target domain performance.

1134 Comparing the TCRL_32_Source curve with the TCRL_48_Source curve, it can be seen that increas-
 1135 ing the total number of target domain environments will reduce the reward value obtained by TCRL
 1136 in the source domain. This means that increasing the total number of target domain environments
 1137 will reduce TCRL’s performance in the source domain. Overall, the default configuration of the total
 1138 number of target domain environments is a more suitable training parameter configuration.

1139 In the Similar Target Environment Number experiment, we change the proportion of environments
 1140 with identical configuration in the total target domain environments, as shown in the right part of
 1141 Fig. 8. The TCRL_4_Source curve and the TCRL_4_Unauth_Target curve represent the default train-
 1142 ing configuration of the TCRL algorithm. In the source domain, different configurations achieve
 1143 similar reward values, with similar trends for the TCRL_4_Source curve, the TCRL_2_Source
 1144 curve and the TCRL_8_Source curve. In the target domain, the reward value obtained by the
 1145 TCRL_2_Unauth_Target curve is higher than that of the other two dotted lines, which means that
 1146 this configuration weakens the performance suppression effect of TCRL on the target domain. That
 1147 is to say, the training performance is poor when the target domain environment where data is col-
 1148 lected by every two threads is the same. Overall, the default configuration of the identical target
 1149 environment proportion in the total target domain environments is a more suitable training param-
 1150 eter configuration.

1151 F.3 REWARD SCALE FOR TARGET DOMAIN ENVIRONMENTS

1152 In this experiment, we mainly aim to verify the impact of target domain environments with different
 1153 random parameters on the final reward value obtained by the policy model. As shown in Fig. 9, the
 1154 final reward value obtained by the policy model in the randomized target domain environment using
 1155 the same randomization control parameters has a significant variance. To ensure the objectivity of
 1156 the experimental results, we scale the reward values based on the final reward values obtained in the
 1157 source and target domains.
 1158



1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170 Figure 9: The impact of target domain environment with different randomization parameters on the
 1171 final reward value. The *Source* line represents the reward curve on the source domain environment,
 1172 and the other three lines *Target_seedx* indicate the reward curves on the different randomized target
 1173 domain environments.

1174
 1175 F.4 HYPERPARAMETER SENSITIVITY

1176 We conducted experiments on the HalfCheetah-v3 environment to analyze the sensitivity of key
 1177 hyperparameters in our framework:

1178 **δ for environment perturbation.** Since δ is determined through iterative optimization as described
 1179 in our Q2 response, we tested how reducing this parameter affects protection capability:
 1180
 1181

1182 Table 4: Sensitivity to environment perturbation δ . “orig.” refers to the original reward.

1183
 1184
 1185
 1186

Perf.	δ (orig.)	$\delta/2$	$\delta/4$	PPO_Trans
HalfCheetah-v3	2516	2558	2617	4115

1187 TCRL maintains effective protection even with reduced perturbation amplitude.

$T_{threshold}$ for quick transfer. This parameter defines the time threshold for identifying quickly transferable environments (20% of training time from scratch). Testing variations of this threshold shows:

Table 5: Sensitivity to quick transfer threshold $T_{threshold}$.

Perf.	$0.8T_{threshold}$	$T_{threshold}$	$1.2T_{threshold}$
Unauth	2461	2516	2671
Auth	4195	4207	4015

The parameter exhibits moderate sensitivity without substantially impacting protection.

MMD loss weight η . This parameter balances feature distribution separation between domains. Testing values around our default (3e-5):

Table 6: Sensitivity to MMD loss weight η .

Perf.	1e-5	3e-5	5e-5
Unauth	2608	2516	2497
Auth	4224	4207	4195

Results show low sensitivity within this range.

KL divergence weight λ . This parameter controls the influence of authorized policy behavior:

Table 7: Sensitivity to KL divergence weight λ .

Perf.	1e-2	1e-3	1e-4
Unauth	2647	2516	2623
Auth	3872	4207	4007

λ shows higher sensitivity than other parameters. Our default value (1e-3) provides optimal balance between maintaining authorized performance while limiting unauthorized performance.

Most parameters show low to moderate sensitivity, with λ requiring the most careful tuning.

F.5 EXPERIMENTAL RESULTS OF OTHER MUJoCo ENVIRONMENTS

The experimental results of the other three MuJoCo environments, such as InvertedDoublePendulum, Walker2d and Humanoid, as shown in Figure 10 and Fig. 11.

In Fig. 10, the baseline PPO algorithm and our TCRL algorithm can achieve similar rewards in the source domain. During the training process, the rewards obtained in the target domain are much less than the rewards in the source domain. In particular, in the Humanoid environment, the green reward curve of the target domain has basically no upward trend.

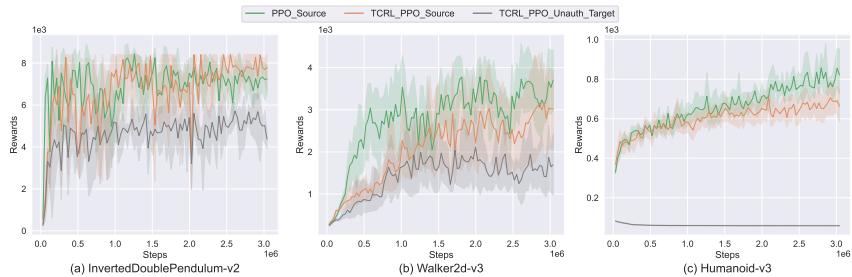
In Fig. 11, the orange reward curves initialized by our TCRL model achieve the worst results, which means that the TCRL model can prevent the migration of the policy model from the source domain to the unauthorized target domain to a certain extent. Meanwhile, the blue reward curves initialized by the original PPO model can obtain similar results with the green reward curves of random initialization. It means that the original PPO model cannot prevent the source domain policy models transfer to the target domain.

In general, these experimental results are similar to those of the three mujoco environment experiments in the main text, which can support the relevant statements in the main text.

F.6 EXPERIMENTAL RESULTS ON HAND MANIPULATION SUITE ENVIRONMENT

In this experiment, we are examining the impact of two transfer reinforcement learning algorithms, namely the DAPG algorithm Rajeswaran et al. (2017a) and the REvolveR algorithm Liu et al. (2022),

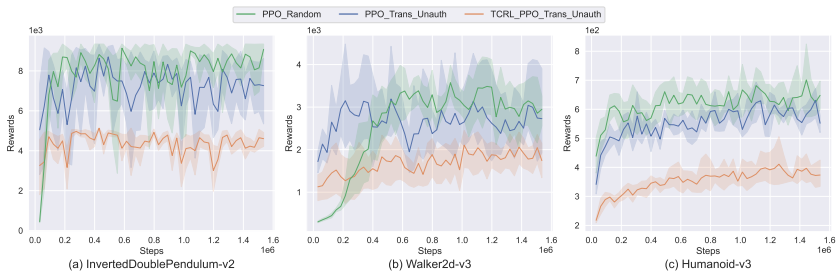
1242
1243
1244
1245
1246
1247
1248
1249
1250



1251
1252
1253
1254
1255

Figure 10: Training performance of the baseline PPO algorithm and our TCRL algorithm. The blue *PPO_Source* and orange *TCRL_PPO_Source* solid curves denote the performance of PPO and TCRL on the source domain, while the green *TCRL_PPO_Unauth_Target* dotted curves indicate the performance of TCRL on the target domains.

1256
1257
1258
1259
1260
1261
1262
1263
1264



1265
1266
1267
1268
1269

Figure 11: Comparing the transfer performance of the PPO and TCRL models on the target domain, the blue *PPO_Trans_Unauth* curve denotes the PPO model and the orange *TCRL_PPO_Trans_Unauth* curve denotes the TCRL model. The blue *PPO_Random* curve, trained with a random initialized model, serves as the baseline.

1270
1271
1272
1273
1274

on the transfer-controllability of the TCRL model. The objective is to evaluate the effectiveness of these algorithms in attacking the transfer-controllability of the TCRL model. In these experiments, we replaced the PPO algorithm in the main text with the NPG algorithm.

1275

F.6.1 HAND MANIPULATION SUITE ENVIRONMENT

1276
1277
1278
1279
1280

This part of the experiment is carried out on the hand manipulation suite environment Liu et al. (2022). This environment is constructed based on the ADROIT platform Rajeswaran et al. (2017a), as shown in Fig. 12.

1281
1282
1283
1284
1285
1286
1287

In Fig. 12, the ADROIT platform is a 24-DoF anthropomorphic platform designed for addressing challenges in dynamic and dexterous manipulation. The first, middle, and ring fingers have 4 degrees of freedom (DoF). Little finger and thumb have 5 DoF, while the wrist has 2 DoF. Each DoF is actuated using position control and is equipped with a joint angle sensor. In this experiment, we use two kinds of these tasks, the object relocation task and the door opening task. As shown in Fig. 12 (a), the goal of the object relocation task is to move the blue ball to the green target. As shown in Fig. 12 (b), the goal of the door opening task is to undo the latch and swing the door open.

1288
1289
1290
1291
1292
1293

The hand manipulation suite environment Liu et al. (2022) is designed to make some evolving transferable environments for transfer reinforcement learning, as shown in Fig. 13. The evolutionary generation process of the transferable five-finger dexterous hand robot is shown in Fig. 13 (c). In the beginning, the hand robot had five dexterous fingers. In the process of continuous evolution, the middle finger, ring finger, and little finger of the robot are getting shorter and shorter. In the end, the hand robot only retained two fingers such as the thumb and index finger, and only had 1 DoF.

1294
1295

Next, we can construct the transferable learning tasks as shown in Fig. 13 (a) and (b). In the object relocation transfer task, the objective of the source domain task is to move a blue ball to the green target using the original five-finger dexterous hand robot. However, in this case, the robot is substi-

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

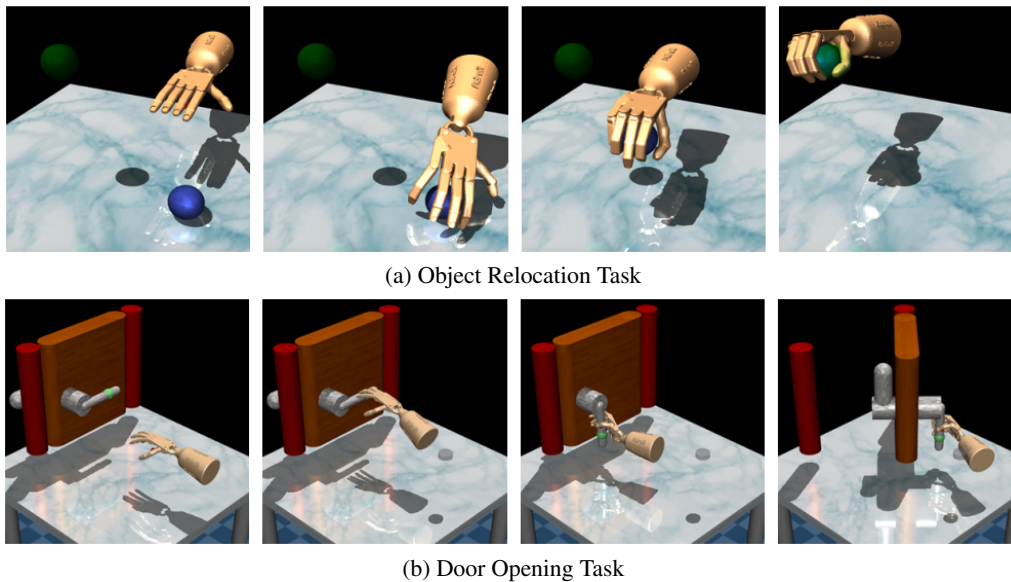


Figure 12: The five-finger dexterous hand provided in the ADROIT platformRajeswaran et al. (2017a).

tuted with a simpler two-finger robot in the target domain. Similarly, in the door opening transfer task, the robot configuration remains unchanged, but the objective is modified to opening the door.

F.6.2 EXPERIMENTAL RESULTS AND ANALYSIS

The experimental results are shown in Table 8 and Table 9 below. In these tables, "From Scratch" means training the policy model from scratch on the target domain, while "Direct Finetune" means using a pre-trained policy model from the source domain for transfer learning. There are two kinds of pre-trained policy model, the "NPG" modelRajeswaran et al. (2017b) and our "TCRL" model. Then, two kinds of transfer reinforcement learning algorithms, the "DAPG" algorithmRajeswaran et al. (2017a) and the "REvolveR" algorithmLiu et al. (2022), are applied to attack the transfer-controllability of the TCRL model. In the "Sparse Reward" setting, only task completion is rewarded. In the "Dense Reward" setting, a distance reward is provided at every step.

In the REvolveR algorithmLiu et al. (2022) and the DAPG algorithmRajeswaran et al. (2017a), an adaptive training scheduling strategy is employed to enhance training efficiency. Consequently, it is not possible to predefine the total number of RL iterations in order to compare performance fairly under the same number of iterations. Instead, the REvolveR algorithmLiu et al. (2022) compares the number of RL optimization steps required to achieve a 90% success rate on the tasks. In this paper, we continue to use the above evaluation method.

From Table 8, none of the transfer learning algorithms initialized with the TCRL model could converge within 100K iterations. The reason may be that the five-finger robot and the two-finger robot grab the blue ball in completely different ways, as shown in Fig. 13 (a). In the TCRL model, due to the reverse training on positive samples in the evolutionary training process, it becomes challenging for transfer reinforcement learning algorithms to obtain positive samples of grasping the blue ball in the target domain. This significantly amplifies the training difficulty for the two-finger robot in the target domain. As a result, the training speed of transfer reinforcement learning using the TCRL model as the initialization model is significantly slowed down in the object relocation task. In other words, the TCRL model has hindered the transfer progress of the DAPG algorithm and the REvolveR algorithm.

From Table 9, the convergence speed of the transfer learning algorithm initialized with the TCRL model is significantly reduced. Compared with the object relocation task, in the door opening task, the execution process of pushing the door handle is similar for the five-fingered robot and the two-

1350
 1351
 1352
 1353
 1354
 1355
 1356
 1357
 1358
 1359
 1360
 1361
 1362
 1363
 1364
 1365
 1366
 1367
 1368
 1369
 1370
 1371
 1372
 1373
 1374
 1375
 1376
 1377
 1378
 1379
 1380
 1381
 1382
 1383
 1384
 1385
 1386
 1387
 1388
 1389
 1390
 1391
 1392
 1393
 1394
 1395
 1396
 1397
 1398
 1399
 1400
 1401
 1402
 1403

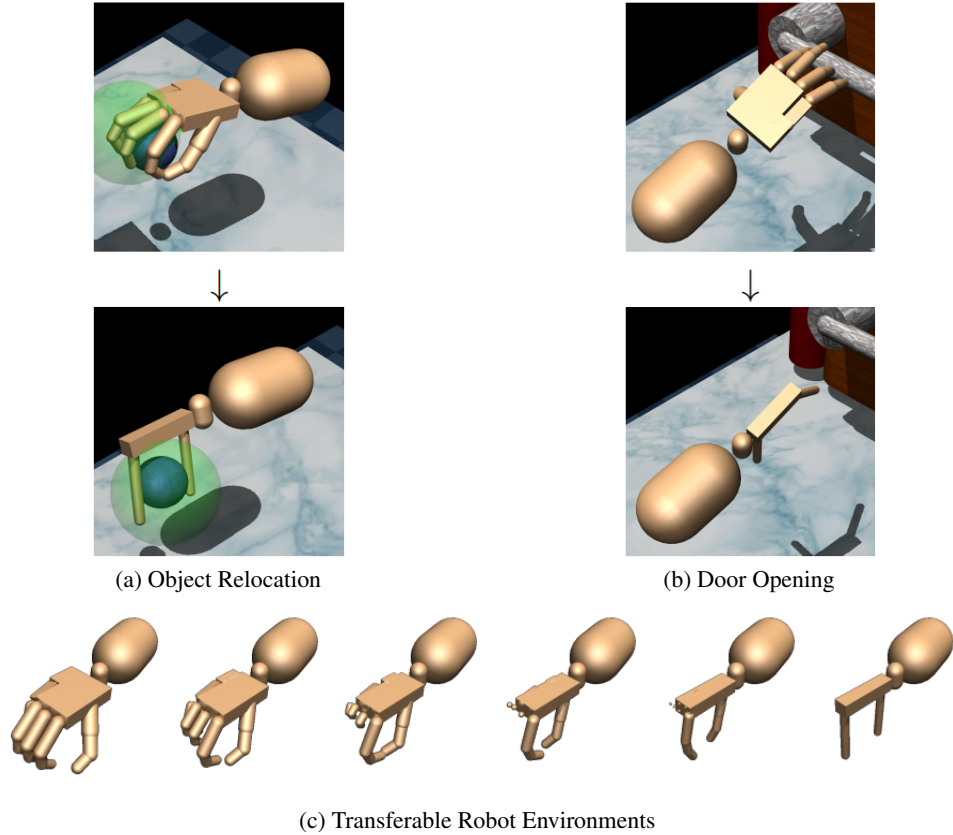


Figure 13: The transferable tasks on hand manipulation suite Liu et al. (2022).

Table 8: The experimental results of the target transfer task

	Dense Reward		Sparse Reward	
	NPG	TCRL	NPG	TCRL
From Scratch	>100K		1	
Initialized Model	NPG	TCRL	NPG	TCRL
Direct Finetune	43.5K	>100K	1	-
DAPGRajeswaran et al. (2017a)	23.3K	>100K	1	-
REvolveRLiu et al. (2022)	-	>100K	18.1K	>100K

1404 fingered robot, as shown in Fig. 13 (b). Therefore, in this task, even with the initialization of the
 1405 TCRL model, the REvolveE algorithm Liu et al. (2022) can still achieve the goal of a success rate
 1406 exceeding 90%. However, our TCRL model can still significantly slow down the convergence speed
 1407 of the REvolveE algorithm, which can still generate certain value in practical applications.

1408
 1409 Table 9: The experimental results of the door opening transfer task

	Dense Reward		Sparse Reward	
From Scratch	-		1	
Initialized Model	NPG	TCRL	NPG	TCRL
Direct Finetune	7.6K	82.5K	1	-
DAPGRajeswaran et al. (2017a)	5.4K	48.3K	1	-
REvolveRLiu et al. (2022)	-	45.4K	2.6K	58.7K

1419 Overall, the above experimental results demonstrate that the TCRL model provides a certain level of
 1420 protection for the intellectual property of the policy model when facing attacks from certain transfer
 1421 reinforcement learning algorithms.

1422 F.7 COMPARISON WITH DOMAIN RANDOMIZATION

1423 While traditional domain randomization (e.g., MAML) aims to enhance generalization, TCRL se-
 1424 lectively restricts transfer to unauthorized domains. Our supplementary experiments demonstrate
 1425 TCRL’s superior performance:

1426
 1427 Table 10: Performance comparison of MAML and TCRL across different domains and environ-
 1428 ments. Values represent rewards.

Method	Domain	HalfCheetah-v3	Hopper-v3	Ant-v3
MAML	Unauthorized	2916	1475	1387
	Authorized	3972	1837	1678
TCRL	Unauthorized	2516	1028	1043
	Authorized	4207	2075	2427

1431
 1432 These results confirm that directly applying domain randomization techniques to our task would lead
 1433 to suboptimal outcomes. Our approach with MMD loss and KL divergence constraints achieves the
 1434 desired balance: limiting performance in unauthorized domains while maintaining or improving it
 1435 in authorized ones.

1442 G DISCUSSION

1443
 1444 **Question1:** To protect the policy model, it is advisable to conceal the model parameters and strictly
 1445 restrict access to an API interface specifically designed for querying policy decisions based on the
 1446 observed state. Given this approach, is it still necessary to implement a transfer-controllable policy?
 1447

1448
 1449 **Answer:** Yes, it is still necessary. Suppose Company A has designed a robot R_A and trained the
 1450 corresponding baseline policy model π_A . At the same time, Company B has replicated a robot R_B
 1451 with similar dynamic characteristics and obtained the API of Company A’s robot’s policy model
 1452 π_A . In this case, Company B can use the API to collect the motion trajectories T_r of robot R_B and
 1453 then use relevant methods of offline reinforcement learning to obtain an approximate version of the
 1454 policy model $\hat{\pi}_A$. By applying transfer learning to the $\hat{\pi}_A$ model, Company B can obtain a suitable
 1455 policy model π_B for robot R_B .

1456 However, when Company A trains the baseline policy model π_A using the TCRL algorithm, if
 1457 Company B tries to use the same API, they would only collect poor-quality motion trajectories. As
 a result, subsequent offline reinforcement learning and transfer learning processes cannot be carried

1458 out. Therefore, training a transfer-controllable policy model becomes necessary in order to mitigate
1459 this issue.

1460
1461 **Question2:** In the paper, the unauthorized target domain environments are designed by randomizing
1462 some parameters in the environments. However, it would be quite rare that the real target application
1463 is only a few parameters different from the source environments while all other settings are the same.

1464 **Answer:** Yes, perhaps such cases are quite rare. However, if Company B intends to steal the intel-
1465 lectual property of Company A's policy model, they would need to take certain steps to construct
1466 a series of similar target domain environments. For example, as shown in Fig. 13 (c), Company B
1467 can create a series of intermediate robots that allow Company A's five-finger hand robot to transition
1468 naturally to Company B's two-finger hand robot. In general, by using transfer learning algorithms,
1469 Company B can avoid some of the errors that Company A would encounter when training from
1470 scratch.

1471 **Question3:** The environment randomization module can be time-consuming and may not be suitable
1472 for all scenarios.

1473 **Answer:** No single method can be universally applicable to all scenarios, and the environment
1474 randomization module is merely a simple preliminary solution. This paper aims to raise awareness
1475 about the issue of protecting policy model intellectual property and propose a general solution. In
1476 practical applications, various more efficient environment randomization schemes can be designed
1477 for this module.

1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511