

# Robots are Exploring to Learn From Each Other: Multi-Robot Multi-Task Transfer

## Abstract

**Motivation:** In this project, we explored a multi-robot, multi-task transfer framework with the objective of scalability across various tasks and robots. This framework is designed to be adaptable to different dynamics and reward structures, emphasizing data efficiency and support for target robots and tasks, especially those with limited resources. This lies in the observation that similar agents often share structural similarities, noise profiles, and challenges. Therefore, enabling robots to learn from each other could prove more effective than starting from scratch. This approach aligns with the initial human fascination with robotics: the aspiration to create machines capable of performing tasks in a human-like manner. An ideal robot would be able to observe human actions and learn directly from them. However, the differing dynamics and morphologies between humans and robots introduce significant noise in action learning and transfer. Understanding how to learn from these differences is crucial, particularly for offline learning and learning from videos. With humans exhibiting varied dynamics, it becomes essential to consider these variations in our learning models. Furthermore, we explore the potential of optimizing robot dynamics. Can we demonstrate that learning from robots is more effective than learning from other sources? A critical question we address is how to effectively train a low-resource robot alongside a high-resource robot in mastering a new task, considering their differing dynamics and morphologies.

**Research Question:** Our primary research challenge involves a scenario where we possess a well-resourced source robot trained on a specific task, contrasted with a low-resource robot characterized by distinct morphology and dynamics. Additionally, we have a different target task, which also has limited resources but differs in objectives and reward structures. Our goal is to propose a transfer learning method that facilitates efficient learning between these disparate entities with minimal interaction. This method should enable the target robot to learn the source task and vice versa. We also explore various approaches to address this problem, aiming to enhance the effectiveness and efficiency of transfer learning in robotics, especially under resource constraints.

**Method:** Our approach employs reward compensation as a key strategy for facilitating dynamics transfer in multi-robot, multi-task scenarios. This same principle of reward compensation is also applied to task transfer. Essentially, our method utilizes reward compensation as a universal tool for both dynamics and task transfers in a multi-robot context. Specifically, we adopt the method outlined in [4] for dynamic transfer, leveraging the differences in rewards as a basis for reward compensation during task transfer.

**Novelty: Distinct Approach:** Our setup introduces a significantly unique approach. Our literature review indicates that

none of the previous works have implemented a setup like ours, distinguishing it from other approaches by emphasizing a more data-efficient robot interaction.

**Reward Compensation for Multi-Task, Multi-Robot Transfer:** We introduce reward compensation as a key strategy in our multi-task, multi-robot transfer framework.

**Enhanced Exploration through Reward Compensation:** We incorporate exploration strategies based on Reward Compensation to improve convergence. This is based on the principle that higher reward compensation indicates a greater need for exploration.

**Mutual Training of Source and Target Agents (Co-Training):** Our method involves co-training source and target agents, utilizing reward compensation as the basis for this mutual learning process.

**Versatile Structure for Efficient Transfer Learning:** We propose various structural models that are supported by reward compensation, aiming to enhance efficiency in multi-robot, multi-task transfer learning.

**Focusing on Low-Resource and High-Resource Robotics:** Our approach takes into account the disparities in resource availability, catering to both low-resource and high-resource robotics scenarios.

**Experiments:** In our experiments, we focus on reacher tasks involving arms with varying morphologies (differing in size and density) in both source and target scenarios. These tasks are further diversified by implementing distinct goals and reward functions for the source and target. Benchmark methods, employing RL algorithms, are applied in both the source and target settings. We have conducted several experiments to assess and analyze the effectiveness of our proposed method through comprehensive ablation studies.

**Results and analysis:** Overall, our analysis of various approaches indicated that both source and target training achieved satisfactory coconvergence. Initially, I anticipated that incorporating an exploration component would lead to faster convergence. However, this exploration actually resulted in greater variability in the rewards, suggesting that less exploration might be more suitable for the current tasks and for more complicated tasks this exploration is more helpful. Nonetheless, the system ultimately approached convergence, implying that with extended training, results could further improve. Additionally, cotraining initially introduced high variability, but this stabilized towards the end, yielding good outcomes. This is because when we engage in the training procedure, we achieve a better estimate of reward compensation, which is a main info sharing in the cotraining.

# 1 Introduction to Approach

The overall approach of the multi-robot multi-task method is depicted in Figure 1. We have a target world, which is a low-resource world, and are interested in training the target task on the target agents using a high-resource world. To achieve this, we are utilizing reward compensation based on action\_reward[4] and task\_reward compensation

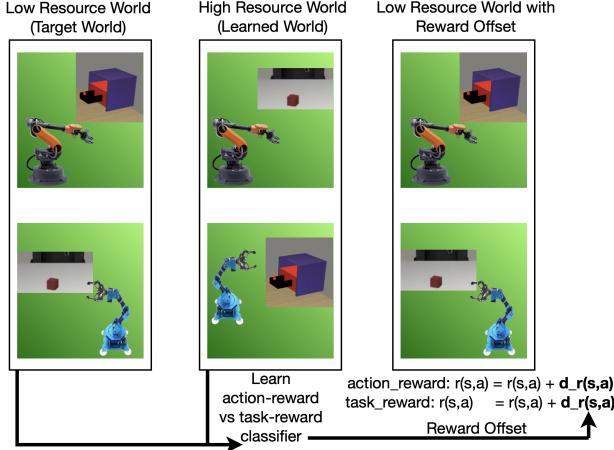


Figure 1: Overall approach

This approach can be extended to multi-task learning for a low\_resource target robot, as shown in Figure 2

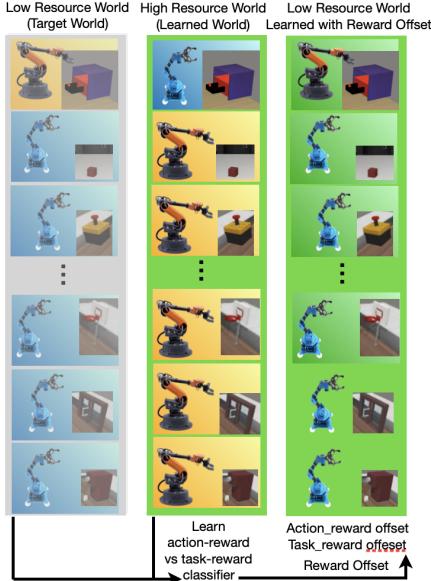


Figure 2: low\_resource robot learning in multi-task learning

Additionally, this Low\_resource transfer learning can be utilized for learning from human demonstrations, as illustrated in Figure 3.

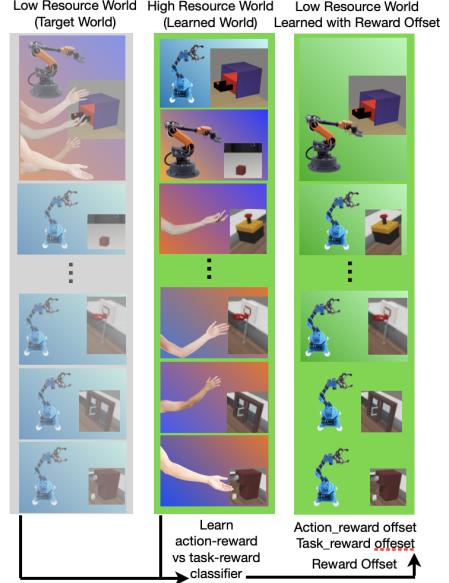


Figure 3: Low\_resource learning from human demonstration.

## 2 Related Work

Bousmalis et al.'s RoboCat [1] is a decision transformer for mastering new skills and embodiments in robotics. Dasari et al.'s 'RoboNet' [2] fuses a substantial dataset from 7 robots with visual foresight and supervised inverse models, showcasing superior adaptability to new tasks and robots. Devin et al. 's method [3] splits neural network policies into 'task-specific' and 'robot-specific' modules, facilitating transfer learning across various robots and tasks. Eysenbach et al. [4] propose a method for domain adaptation in reinforcement learning, which uses a modified reward function learned from auxiliary classifiers and has been proven to recover a near-optimal policy for the target domain. Open X-Embodiment [5] introduces RT-X, a high-capacity model trained on data from 22 robots, enhancing multiple robots' capabilities through positive transfer and mutual learning, marking a step towards seamless task adaptation in robotics.

## 3 Methodology

### 3.1 Problem Statement

Robot  $R_1$  has learned policy  $\pi_1$  for Task  $T_1$ , and robot  $R_2$  has learned policy  $\pi_2$  for Task  $T_2$ .  $R_1$  and  $R_2$  have different dynamics, and  $T_1$  and  $T_2$  have different reward functions. Our goal is for  $R_1$  to learn policy  $T_2$  with a few interactions and for  $R_2$  to learn  $T_1$  with a few interactions (zero-shot, one-shot, few-shot) as well. (Since  $R_2$  knows  $T_2$  and  $R_1$  knows  $T_1$ , this task is an extension of in-context learning.) Another suggested transfer approach is illustrated in Figure 5, and a different approach is depicted in Figure 6. The previous approach requires another framework, Online Robot Learning

Transfer, in a simpler case, as shown in Figure 4. These three methods can occur because of separate action and task transfer by reward compensation.

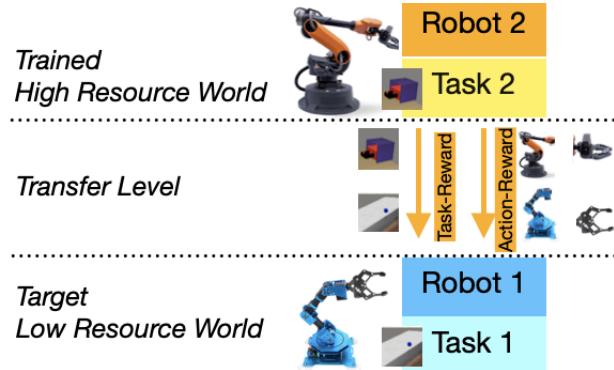


Figure 4: Single Transfer: Two robots learning from each other: Transfer learning for multi-robot, multi-task scenarios involving two robots and two tasks.

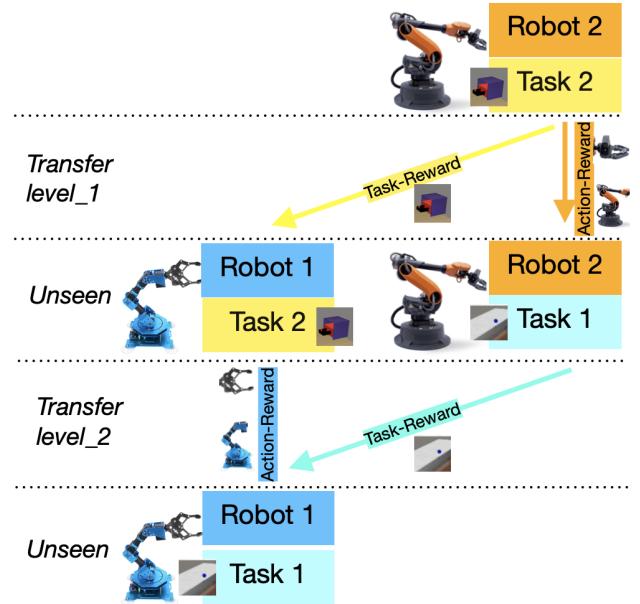


Figure 6: Type\_2 of Transfer

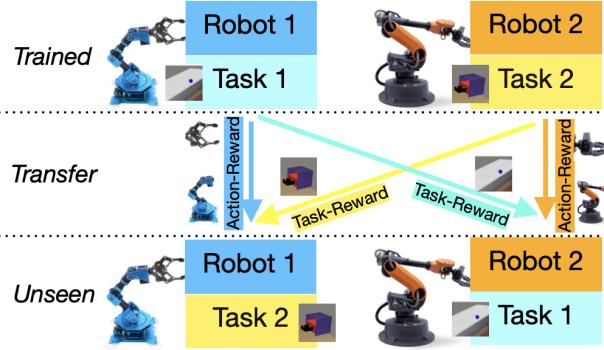


Figure 5: Type\_1 of Transfer

### 3.2 Method overview

You can check the method overview for the setup of transfer1 in Figure 7. A similar approach applies to transfer2 as well, which is explained in the algorithm.

### 3.3 Possibility and Analysis

The initial reason that why the methodology work is that is we can use reward editing for having a same dynamics but different tasks (task transfer), and also we can use reward editing for learning a same task for different dynamics (dynamics transfer). So, apparently we can can use reward editing when we have two different dynamics and each know different tasks. and these two can learn from each other. Simple sequence of this transfer can help. The point is that how to use efficient exploration to have the efficient training of the agent in online one. The point is that we limit the exploration based on the reward editing formulation that we learned from  $\pi_{1,1}$  and  $\pi_{2,2}$ . Because we believe that the state\_action space

is well explored with the  $\pi_{1,1}$  and  $\pi_{2,2}$  (we can have a random exploration on the other side as well for the state\_action that never seen before and based on how much it improves the result, changing the  $\theta$  which shows how much out of state\_action exploration needed).

### 3.4 Algorithm

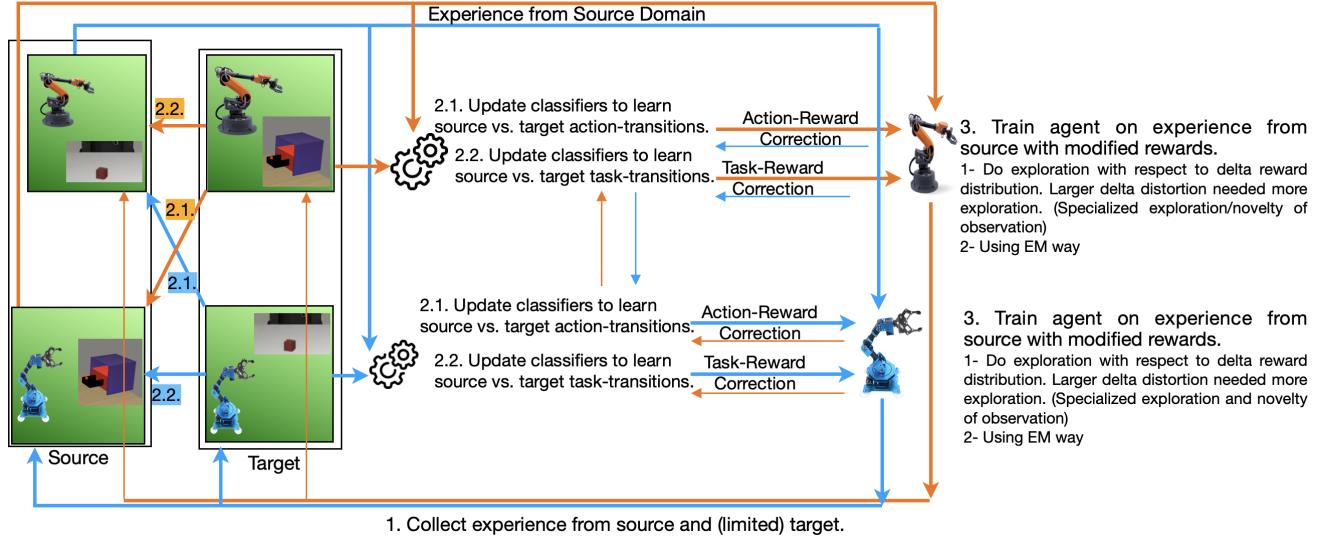
We aim to cover three types of transfer for multi-robot multi-agent systems, as depicted in Figures 4, 5, and 6, using reward compensation. For dynamics transfer, we are utilizing reward compensation from [4] as action\_reward compensation, and for task\_reward, we are using  $\delta_r$  for each action\_pair separately in each transfer.

Subsequently, we can introduce the following algorithms for multi\_robot multi\_task transfer: 1, 2, and 3. Considering that reward compensation indicates how far the target action state space is from the source action state space, we use reward compensation as an exploration parameter. Moreover, training is conducted based on the reward compensation learned from the target to update the source policy.

## 4 Experiments and Results

The experimental setup is as follows: The task involves a reacher task featuring different reward functions, varying robot dynamics and morphology, as well as different densities. It includes two distinct arms, each with unique dynamics, goal locations, and reward formulas. One arm uses the Euclidean distance for reward calculation, while the other employs the Manhattan distance.

To examine various ablation studies related to our distinct elements of reward compensation in the context of




---

**Algorithm 1** Multi-Robot/Multi-Task(Multi-Reward) Transfer

---

```

1: Input: source MDP  $M_{\text{source}}$ , target MDP  $M_{\text{target}}$ , ratio  $r$  of experience from source vs. target
2: Initialize: replay buffers  $D_{\text{source}}, D_{\text{target}}$ , policy  $\pi$ , parameters  $\theta = (\theta_{\text{SAS}}, \theta_{\text{SA}})$ 
3: for  $t = 1, \dots, \text{num iterations}$  do
4:    $D_{\text{source}} \leftarrow D_{\text{source}} \cup \text{ROLLOUT}(\pi, M_{\text{source}})$  ▷ Collect source data
5:   if  $t \bmod r = 0$  then
6:      $D_{\text{target}} \leftarrow D_{\text{target}} \cup \text{ROLLOUT}(\pi, M_{\text{target}})$  ▷ Collect target data
7:   end if
8:    $\theta \leftarrow \theta - \eta \nabla_{\theta} l(\theta)$  ▷ Update parameters
9:    $\tilde{r}(s_t, a_t, s_{t+1}) \leftarrow r(s_t, a_t) + \Delta r(s_t, a_t, s_{t+1})$  ▷ Compute reward compensation (action-reward transfer or task-reward transfer)
10:   $\pi \leftarrow \text{MAXENT RL}(\pi, D_{\text{source}}, \tilde{r})$  ▷ Update policy
11: end for
12: return  $\pi$ 

```

---

transfer learning, we consider the following setup for the experiments: (In all the experiments, RL on the source and target is considered as the benchmark for each target or source policies, applied separately to the source and target worlds and computed and illustrated in all the diagrams for reference.).

- V1:  $Arm_1$  as a source agent and  $Arm_2$  as a target agent, considering the same reward function as  $Arm_1$ . (Different dynamics, same reward function)
- V2:  $Arm_1$  as a source agent and  $Arm_1$  as a target agent, considering the same reward function as  $Arm_2$ . (Different dynamics, same reward function) For these V1 and V2 experiments, refer to plot 4.
- V3: The dynamics of the source and target are the same as in  $Arm_1$ , but the reward functions differ; one is Euclidean and the other is Manhattan.

- V4: The dynamics of the source and target are the same as in  $Arm_2$ , but the reward functions differ (one is Euclidean and the other is Manhattan). For these V3 and V4 experiments, refer to plot 9.
- V5 – a: In this experiment,  $Arm_1$  runs on  $Task_1$  and transfers to  $Arm_2$  on  $Task_2$ , using the approach of Algorithm 1.
- V5 – b:  $Arm_2$  on  $Task_2$  transfers to  $Arm_1$  on  $Task_1$ , using the approach of Algorithm 1. For these V5 – a and V5 – b experiments, refer to plot 10.
- V6 (reward compensation based exploration with SAC training): Using  $\delta r$  with reward compensation exploration. For the V6 experiments, refer to plot 12.
- V7: Experiments using the  $Transfer_1$  approach and co-training.

---

**Algorithm 2** Single Multi-Robot Multi-Task Transfer

---

```

1: Input: source MDP  $M_{\text{source}}$ , target MDP  $M_{\text{target}}$ , ratio  $r_{\text{target}}$  of experience from source vs. target, ratio  $r_{\text{co}}$  of source training, ratio  $r_{\text{co}}$  of source target co-training
2: Initialize: replay buffers  $D_{\text{source}}, D_{\text{target}}$ , policy  $\pi_{\text{source}}$ , policy  $\pi_{\text{target}} = \pi_{\text{source}}$ ,
3: parameters  $\theta_{\text{st}} = ((\theta_{\text{SAS\_action/dynamics}}, \theta_{\text{SA\_action/dynamics}}), (\theta_{\text{SAS\_task}}, \theta_{\text{SA\_task}}))$ 
4: parameters  $\theta_{\text{ts}} = ((\theta_{\text{SAS\_action/dynamics}}, \theta_{\text{SA\_action/dynamics}}), (\theta_{\text{SAS\_task}}, \theta_{\text{SA\_task}}))$ 
5: for  $t = 1, \dots, \text{num\_iterations}$  do
6:    $D_{\text{source}} \leftarrow D_{\text{source}} \cup \text{ROLLOUT}(\pi, M_{\text{source}})$  ▷ Collect source data
7:   if  $t \bmod r_{\text{target}} = 0$  then
8:      $D_{\text{target}} \leftarrow D_{\text{target}} \cup \text{ROLLOUT}(\pi, M_{\text{target}})$  ▷ Collect target data
9:      $\theta_{\text{st}} \leftarrow \theta_{\text{st}} - \eta \nabla_{\theta_{\text{st}}} l(\theta_{\text{st}})$  ▷ Update parameters
10:     $\tilde{r}(s_t, a_t, s_{t+1}) \leftarrow r(s_t, a_t) + \alpha \cdot \Delta r_{\text{action/dynamics}}(s_t, a_t, s_{t+1}) + (1 - \alpha) \cdot \Delta \tilde{r}_{\text{action/dynamics}}(s_t, a_t, s_{t+1})$  ▷ Compute reward
11:     $\pi_{\text{target}} \leftarrow \text{MAXENT RL}(\pi_{\text{target}}, D_{\text{target}}, \tilde{r}, \Delta r_{\text{task}}, \Delta r_{\text{action}})$  ▷ Update policy by  $\Delta r$  exploration in entropy
12:   end if
13:    $D_{\text{source}} \leftarrow D_{\text{source}} \cup \text{ROLLOUT}(\pi, M_{\text{source}})$  ▷ Collect source data
14:    $\theta_{\text{ts}} \leftarrow \theta_{\text{ts}} - \eta \nabla_{\theta_{\text{ts}}} l(\theta_{\text{ts}})$  ▷ Update parameters
15:   if  $t \bmod r_{\text{co}} = 0$  then
16:      $\tilde{r}(s_t, a_t, s_{t+1}) \leftarrow r(s_t, a_t) + \beta \cdot \Delta r_{\text{action/dynamics}}(s_t, a_t, s_{t+1}) + (1 - \beta) \cdot \Delta \tilde{r}_{\text{action/dynamics}}(s_t, a_t, s_{t+1})$  ▷ Update policy by training on reward compensation based on the source
17:   else
18:      $\tilde{r}(s_t, a_t, s_{t+1}) \leftarrow r(s_t, a_t) + \Delta r_{\text{action/dynamics}}(s_t, a_t, s_{t+1})$  ▷ Update policy by training on reward compensation based on the source
19:   end if
20:    $\pi_{\text{source}} \leftarrow \text{MAXENT RL}(\pi_{\text{source}}, D_{\text{source}}, \tilde{r}, \Delta r_{\text{task}}, \Delta r_{\text{action}})$  ▷ Update policy by  $\Delta r$  exploration in entropy
21: end for
22: return  $\pi$ 

```

---

- V8: Experiments using the  $\text{Transfer}_2$  approach and co-training. For the V7 and V8 experiments, refer to plot 13.
- V9: Experiments using the  $\text{Transfer}_1$  approach and co-training with reward compensation exploration. and V10: Experiments using the  $\text{Transfer}_2$  approach and co-training with reward compensation exploration. For the V9 and V10 experiments, refer to plot 14.

Also, here is the explanation of the diagram's legends for plots: In all diagram legends, the setup vi is explained as above. The tag  $_{\text{low\_r}}$  indicates that the agent is in a low-resource case, meaning that there are limited interactions. The  $_{\text{high\_r}}$  case implies that the target agent is considered to be in a high-resource scenario, allowing for one-to-one updates in the transfer learning. The tag  $_{\text{algo}}$  signifies that the results are from running the main algorithm of the study in the corresponding experiments. The tag  $_{\text{target}}$  indicates that, instead of using the algorithm of the study, the Reinforcement Learning (RL) algorithm (in this case, SAC) is run as a benchmark on the target, using reward compensation for exploration in transformation (if the source and target are different). The tag  $_{\text{source}}$  implies that, similarly, the RL algorithm (SAC) is run as a benchmark on the target, but with the focus on reward compensation for exploration in transformation (if the source and target differ).

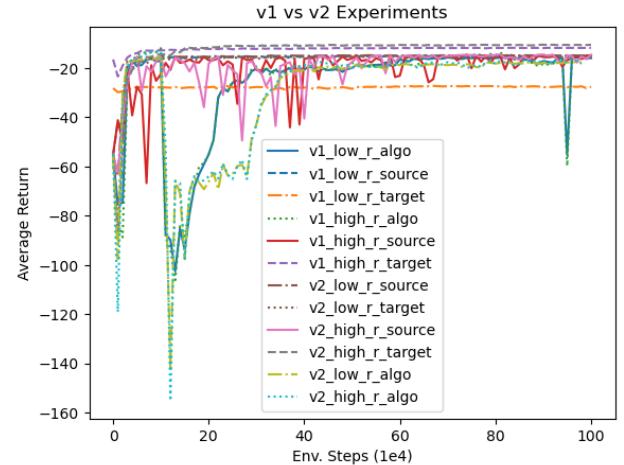


Figure 8: V1:  $Arm_1$  as a source agent and  $Arm_2$  as a target agent, considering the same reward function as  $Arm_1$ . V2:  $Arm_1$  as a source agent and  $Arm_1$  as a target agent, considering the same reward function as  $Arm_2$ . (Different dynamics, same reward function)

## 5 Discussion and Analysis

- In the V1 and V2 experiments, as shown in Plot 4 ('Different Dynamics, Same Reward Function'), it is evident that  $v1_{\text{low\_r\_algo}}$  and  $v1_{\text{high\_r\_algo}}$  exhibit sim-

---

**Algorithm 3** Multiple Multi-Robot Multi-Task Transfer

---

```

1: Input: source MDP  $M_{11}$  and  $M_{22}$ , target MDP  $M_{12}$  and  $M_{21}$ , ratio  $r_{target}$  of experience from source vs. target, ratio  $r_{co}$  of source training, ratio  $r_{co}$  of source target co_training
2: Initialize: replay buffers  $D_{22}, D_{11}, D_{12}, D_{21}$  policy  $\pi_{11}, \pi_{22}, \pi_{12}, \pi_{21}$ .
3: parameters  $\theta_{22}, \theta_{11}, \theta_{21}, \theta_{12}$ 
4: parameters  $\theta_{21} = ((\theta_{SAS.action/dynamics}, \theta_{SA.action/dynamics})_{source2}, (\theta_{SAS.task}, \theta_{SA.task})_{source1})$ 
5: parameters  $\theta_{12} = ((\theta_{SAS.action/dynamics}, \theta_{SA.action/dynamics})_{source1}, (\theta_{SAS.task}, \theta_{SA.task})_{source2})$  ▷ The parameters for action_reward and task_reward could be shared across transfers of the same type.
6: for  $t = 1, \dots, \text{num\_iterations}$  do
7:    $D_{11} \leftarrow D_{11} \cup \text{ROLLOUT}(\pi, M_{11})$  ▷ Collect source data
8:    $D_{22} \leftarrow D_{22} \cup \text{ROLLOUT}(\pi, M_{22})$  ▷ Collect source data
9:   if  $t \bmod r_{11} = 0$  then
10:     $D_{12} \leftarrow D_{12} \cup \text{ROLLOUT}(\pi, M_{12})$  ▷ Collect target data
11:     $D_{21} \leftarrow D_{21} \cup \text{ROLLOUT}(\pi, M_{21})$  ▷ Collect target data
12:     $\theta_{11} \leftarrow \theta_{11} - \eta \nabla_{\theta_{11}} l(\theta_{11})$  ▷ Update parameters
13:     $\theta_{12} \leftarrow \theta_{12} - \eta \nabla_{\theta_{12}} l(\theta_{12})$  ▷ Update parameters
14:     $\theta_{22} \leftarrow \theta_{22} - \eta \nabla_{\theta_{22}} l(\theta_{22})$  ▷ Update parameters
15:     $\theta_{21} \leftarrow \theta_{21} - \eta \nabla_{\theta_{21}} l(\theta_{21})$  ▷ Update parameters
16:     $\tilde{r}(s_t, a_t, s_{t+1}) \leftarrow r(s_t, a_t) + \alpha \cdot \Delta r_{action/dynamics}(s_t, a_t, s_{t+1}) + (1 - \alpha) \cdot \Delta \tilde{r}_{action/dynamics}(s_t, a_t, s_{t+1})$  ▷ Compute reward compensation of each policy based on  $\theta$  values of each policy
17:     $\pi_{target} \leftarrow \text{MAXENT RL}(\pi_{target}, D_{target}, \tilde{r}, \Delta r_{task}, \Delta r_{action})$  ▷ Update policy by  $\Delta r$  exploration in entropy
18:   end if
19:    $D_{11} \leftarrow D_{11} \cup \text{ROLLOUT}(\pi, M_{11})$  ▷ Collect source data
20:    $\theta_{11} \leftarrow \theta_{11} - \eta \nabla_{\theta_{11}} l(\theta_{11})$  ▷ Update parameters
21:    $\theta_{22} \leftarrow \theta_{22} - \eta \nabla_{\theta_{22}} l(\theta_{22})$  ▷ Update parameters
22:   if  $t \bmod r_{co} = 0$  then
23:      $\tilde{r}(s_t, a_t, s_{t+1}) \leftarrow r(s_t, a_t) + \beta \cdot \Delta r_{action/dynamics}(s_t, a_t, s_{t+1}) + (1 - \beta) \cdot \Delta \tilde{r}_{action/dynamics}(s_t, a_t, s_{t+1})$  ▷ Update policy by training on reward compensation based on the source
24:   else
25:      $\tilde{r}(s_t, a_t, s_{t+1}) \leftarrow r(s_t, a_t) + \Delta r_{action/dynamics}(s_t, a_t, s_{t+1})$  ▷ Update policy by training on reward compensation based on the source
26:   end if
27:    $\pi_{11} \leftarrow \text{MAXENT RL}(\pi_{11}, D_{11}, \tilde{r}, \Delta r_{task}, \Delta r_{action})$  ▷ Update policy by  $\Delta r$  exploration in entropy
28:    $\pi_{22} \leftarrow \text{MAXENT RL}(\pi_{22}, D_{22}, \tilde{r}, \Delta r_{task}, \Delta r_{action})$  ▷ Update policy by  $\Delta r$  exploration in entropy
31: end for
32: return  $\pi$ 

```

---

ilar learning curves. This similarity suggests that the transfer learning is updating both trainings in a comparable manner. A similar behavior is observed in the V2 version as well. Therefore, in the default case where we have the same reward function but different dynamics, the main approach is deemed reliable, particularly when considering Algorithm 1, which is a simpler setup compared to Algorithms 2 and 3.

- Figure 9 illustrates the transfer case when there are the same dynamics but different reward functions (one is Euclidean and the other is Manhattan, different tasks and goals). As observed, in v3\_low\_r\_algo and v3\_high\_r\_algo, when utilizing Algorithm 1, both low and high resources behave similarly and converge in the same way during transfer learning. This pattern is consistent for V4 as well, indicating that the reward compensation is effective in scenarios where the dynamics are the same but the rewards differ.

- Therefore, we have already determined that in the default case of identical reward and dynamics, this type of reward compensation works in reward compensation transfer learning. Next, let's examine the results of other types of transfer by reward compensation.

- In the configuration where we have different rewards and dynamics (fig. 4) 10, by comparing v5\_low\_r\_algo with v5\_high\_r\_algo and v5\_b\_low\_r\_algo with v5\_b\_high\_r\_algo, it is observed that the low\_resource agent learns at approximately the same speed as the high\_resource agent. Although there are more fluctuations than in the two previous learning curves, the training process follows a similar pattern. Additionally, the return convergence is lower compared to the previous default setup. Let's examine whether adding more exploration to the SAC algorithm using reward compensation aids in achieving faster convergence.

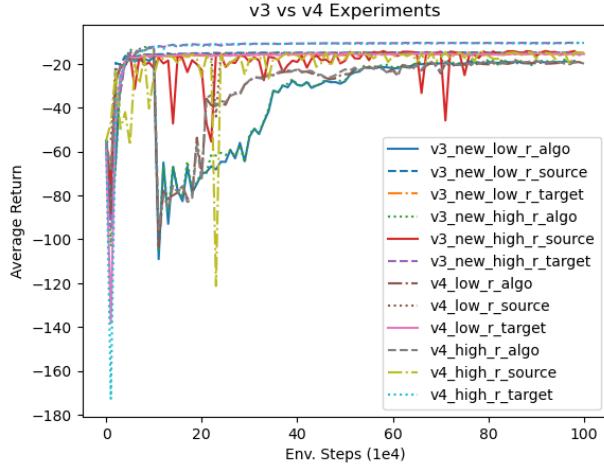


Figure 9: The dynamics of the source and target are the same as in  $Arm_1$ , but the reward functions differ; one is Euclidean and the other is Manhattan. The dynamics of the source and target are the same as in  $Arm_2$ , but the reward functions differ (one is Euclidean and the other is Manhattan).

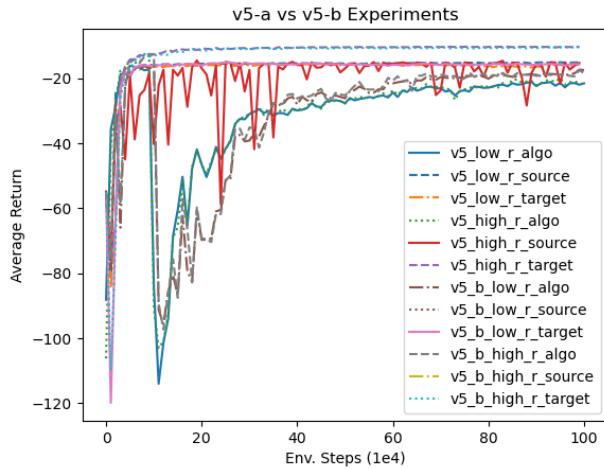


Figure 10: V5 – a: In this experiment,  $Arm_1$  runs on  $Task_1$  and transfers to  $Arm_2$  on  $Task_2$ , using the approach of Algorithm 1. V5 – b:  $Arm_2$  on  $Task_2$  transfers to  $Arm_1$  on  $Task_1$ , using the approach of Algorithm 1.

- Plot 11 illustrates the impact of exploration. In the current setup of the experiment, adding reward exploration appears to be unhelpful for the setup of Experiment 5, introducing more fluctuations in the transfer learning process. We anticipated that sac\_low\_r.algo would converge faster than v5\_b\_low\_r.algo, and sac\_high\_r.algo would converge faster than v5\_b\_high\_r.algo, but neither occurred. Instead, we observed lower rewards and more fluctuations (which are due to increased exploration in the new SAC augmented with reward compensation exploration). However, the exploration-based method eventually achieves the return value of v5. More environment steps may result in a better final return

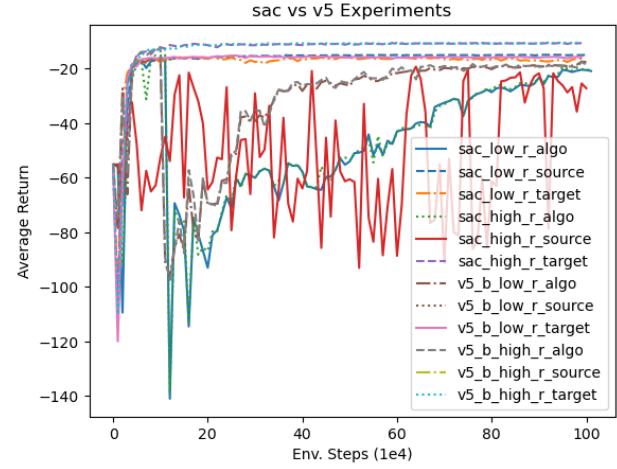


Figure 11: Using  $\delta r$  as reward compensation exploration in SAC algorithm of approach V5 and comparison with approach V5

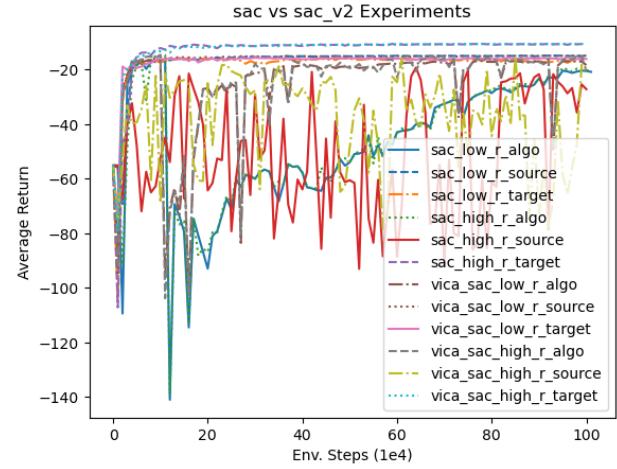


Figure 12: V6: Experiments using the SAC and reward compensation exploration in setup of approach in V5. In tag "vica" we reversed that source and target agents and tasks and redo the experiments with current setup.

value.

- Let's consider changing the target and source, as one of the tasks or agents might be more challenging than the other, and exploration could fail when we use transfer. If this is the case, we should see better results when we swap the target and source. This is because if the previous target's exploration was more complex, then now the source would cover the required state space, making learning targets easier. However, if we observe the same pattern in exploration, it indicates that a more complicated state space is needed for exploration to be effective. Therefore, as shown in 11, we notice better convergence, suggesting that reward compensation exploration was helpful. But we did not get better results by just

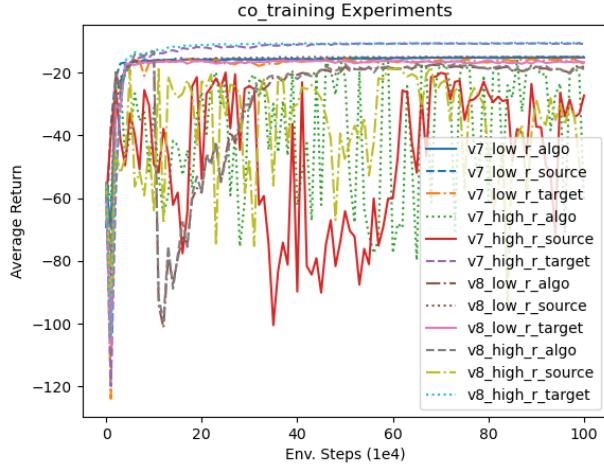


Figure 13: *V7*: Experiments using the *Transfer<sub>1</sub>* approach and co-training in algorithm 2 and *V8*: Experiments using the *Transfer<sub>2</sub>* approach and co-training in algorithm 3

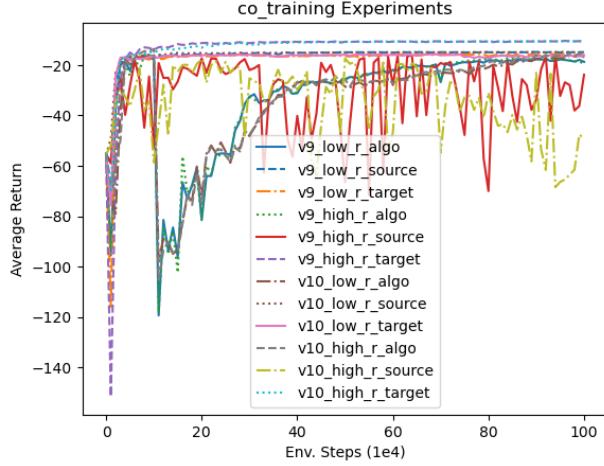


Figure 14: *V9*: Experiments using the *Transfer<sub>1</sub>* approach and co-training in algorithm 2 with reward compensation exploration and *V10*: Experiments using the *Transfer<sub>2</sub>* approach and co-training in algorithm 3 with reward compensation exploration

simply running RL on the target.

- Now, we aim to study the result of co-training on transfer method one and transfer method two in Algorithms 2 and 3. Based on 13, in *V7*, which demonstrates the transfer\_1 method, a comparison between v7\_low\_r\_algo and v7\_high\_r\_algo indicates that there is not effective transfer in co-training, at least in the early stages. However, for v8\_low\_r\_algo and v8\_high\_r\_algo, it is observed that there is good convergence between the source and target in the transfer\_2 method. See 13 for details.
- Plot 14 demonstrates how reward compensation exploration aids in co-training. As observed in *V9* and *V10*, the fluctuations in transfer\_1 become smaller,

indicating that exploration has assisted co-training in transfer\_1 (as seen by comparing v9\_low\_r\_algo with v9\_high\_r\_algo). However, it did not significantly affect the performance of transfer\_2, as shown by the comparison of v10\_low\_r\_algo with v10\_high\_r\_algo.

## 6 Limitations

The main limitation of our project is that the state spaces of the tasks remain identical across different transfer problems. This issue can be mitigated by adopting an observation pixel-based state space. However, a more diverse dataset may be required, depending on the type of task.

## 7 Conclusion and Future works

In this project, we investigated the problem of multi-agent (with different dynamics) and multi-task (with different rewards) transfer. We conducted studies and experiments on reward compensation and exploration based on this compensation in multi-agent and multi-task transformers, using various setups. The types of task transfers can be extended to more agents and more complicated tasks with additional steps. It is likely that a more diverse range of tasks could lead to improved performance. The offline extension of this method is also beneficial. In addition to various types of transfers, the S4 and S5 extensions demonstrate strong performance in in-context RL. I believe that extending reward compensation aids in training sequence models through multi-task, multi-agent transitions and enhances the capacity of different transfer models for multi-agent, multi-task transitions.

## 8 Acknowledgment

The initial idea for this paper originated during the Deep Reinforcement Learning course in the fall of 2023 at UC Berkeley. I would like to thank the instructor team for their support.

## References

- [1] Bousmalis, Konstantinos, Giulia Vezzani, Dushyant Rao, Coline Devin, Alex X. Lee, Maria Bauza, Todor Davchev et al. "RoboCat: A Self-Improving Foundation Agent for Robotic Manipulation." arXiv preprint arXiv:2306.11706 (2023).
- [2] Dasari, Sudeep, Frederik Ebert, Stephen Tian, Suraj Nair, Bernadette Bucher, Karl Schmeckpeper, Siddharth Singh, Sergey Levine, and Chelsea Finn. "Robonet: Large-scale multi-robot learning." arXiv preprint arXiv:1910.11215 (2019).
- [3] Devin, C., Gupta, A., Darrell, T., Abbeel, P., & Levine, S. (2017, May). Learning modular neural network policies for multi-task and multi-robot transfer. In 2017 IEEE international conference on robotics and automation (ICRA) (pp. 2169-2176). IEEE.
- [4] Eysenbach, B., Asawa, S., Chaudhari, S., Levine, S., & Salakhutdinov, R. (2020). Off-dynamics reinforcement learning: Training for transfer with domain classifiers. arXiv preprint arXiv:2006.13916.

[5] Padalkar, Abhishek, Acorn Pooley, Ajinkya Jain, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky et al. "Open x-embodiment: Robotic learning datasets and rt-x models." arXiv preprint arXiv:2310.08864 (2023).