# CLIP-RL: Aligning Language and Policy Representations for Task Transfer in Reinforcement Learning

**Chainesh Gautam**
Department of Data Science and Artificial Intelligence
International Institute of Information Technology
Bangalore, IN 560100
chainesh.gautam@iiitb.ac.in

**Raghuram Bharadwaj Diddigi**
Department of Data Science and Artificial Intelligence
International Institute of Information Technology
Bangalore, IN 560100
raghuram.bharadwaj@iiitb.ac.in

## Abstract

Deep Reinforcement Learning (RL) provides algorithms for solving complex sequential decision-making problems by leveraging the function approximation capabilities of neural networks. Recently, there has been an increasing need to develop agents capable of solving multiple tasks within the same environment, especially when these tasks are naturally associated with language. For instance, in a warehouse setting, a robot may be tasked with executing instructions like, "Go to location A, pick up object B, and drop it at location C." In such scenarios, it would be beneficial to transfer knowledge from previously trained similar tasks to new ones, rather than training models from scratch for each new task. A straightforward approach in this setting might involve extracting sentence embeddings of natural language instructions, identifying the closest pre-trained instruction using similarity metrics (such as cosine similarity), and initializing the new task's neural network with the trained weights of the closest pre-trained task. However, this method is limited by the assumption that language similarity implies policy similarity, which is not always valid. In this work, we propose a novel approach that leverages combinations of pre-trained (language, policy) pairs to establish an efficient transfer pipeline. Our algorithm is inspired by the principles of Contrastive Language-Image Pretraining (CLIP) in Computer Vision, which aligns representations across different modalities under the philosophy that "two modalities representing the same concept should have similar representations." The central idea here is that the instruction and corresponding policy of a task represent the same concept, the task itself, in two different modalities. Therefore, by extending the idea of CLIP to RL, our method creates a unified representation space for natural language and policy embeddings. Experimental results demonstrate the utility of our algorithm in achieving faster transfer across tasks.

# 1    Introduction

We consider a setting where an autonomous agent is tasked with learning a set of tasks using a deep reinforcement learning (RL) algorithm. Each task is succinctly described by a natural language instruction, providing a high-level specification of the goal. A key question in this context is: "When training on a new task, can we leverage knowledge from a subset of tasks the agent has already been trained on?" Efficient transfer is critical in such scenarios, as it can significantly reduce training time and resources. One popular approach to enable transfer is to initialize the neural network for the new task intelligently. A simple method is to identify the pre-trained task whose natural language instruction is most similar to the new task's instruction (using a similarity metric like cosine similarity) and initialize the network with the corresponding learned weights. However, this approach inherently assumes that natural language similarity directly correlates with policy similarity, an assumption that is not necessarily true in practice.

For example, consider a simple grid-world environment illustrated in Figure 1. This environment consists of six tasks where the agent's objective is to navigate to specific locations based on objects identified by a combination of color and shape. Assume the agent has been trained on five tasks, and the sixth task is "go to the red cone." The natural language embeddings for this new instruction might be similar to those for instructions like "go to the red box" or "go to the blue cone" based on standard embedding similarity metrics. However, the structure of the environment organizes objects by color rather than shape, making it more appropriate to prioritize policy similarity derived from color-based navigation. Therefore, the challenge lies in tuning instruction embeddings to reflect the structural properties of the environment, ensuring that transfer aligns with task-specific dynamics rather than superficial language similarity. Notice that the weights of the neural networks corresponding to tasks grouped by color, rather than shape, are more likely to exhibit similar patterns. This is because tasks that share structural properties within the environment, such as navigating based on color, lead to similar policies.
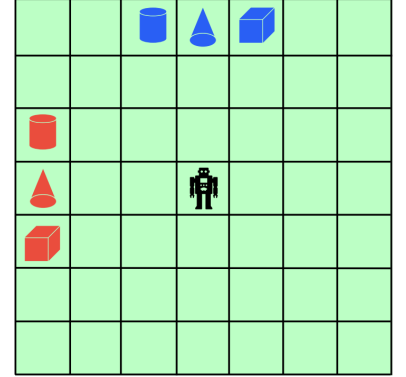


Figure 1: An illustration of grid-world

Therefore, an intelligent transfer approach should leverage this inherent similarity when tuning instruction embeddings. Specifically, the textual description of a task and the neural network weights representing the corresponding policy embody the same underlying concept but in different modalities: language and policy representations, respectively. This dual-modality representation naturally lends itself to a Contrastive Language-Image Pretraining (CLIP) [8] inspired framework. In CLIP, text and image embeddings are aligned so that pairs referring to the same concept are closer in representation space. Analogously, we propose aligning natural language embeddings of task instructions with the weights of the corresponding policy networks. By applying a contrastive learning objective, instruction embeddings are tuned to reflect task-specific policy similarities rather than superficial linguistic resemblance.

As a result, after tuning, the embeddings of instructions capture the structure of the environment, ensuring that tasks with similar policies are appropriately aligned, even if their language descriptions are not closely related. This alignment significantly improves transfer performance, reducing reliance on language-based similarity and enabling more effective task generalization.

To demonstrate the effectiveness of our approach, we conduct experiments on grid-world environments of varying and increasing sizes. Our results show that the CLIP-inspired transfer model reduces training time by approximately 50% on average compared to a baseline model that relies solely on natural language similarity for task transfer. Additionally, we observe that as the grid size increases, the performance gains from our method become even more pronounced. This highlights the scalability and robustness of our approach, particularly in more complex environments where naive language-based transfer methods are less reliable.

# 2    Related Works

Since the time immemorial, humans make use of instruction in more natural form, namely, language to solve difficult problems using the knowledge of relatively similar problem. Meanwhile, teaching agents how to solve a new task is a central problem in Reinforcement Learning. One of the approach to solve this problem is using natural language instruction to design rewards [6, 4] , but as the complexity of the tasks grow it's challenging to give intricate details of the task using natural language. Another paradigm, called imitation learning [1], requires demonstration(s) of the similar desired task, which can then be used to learn the policy for the desired task. However, for each new task, getting the set of demonstration requires lot of human efforts and is not scalable.

In the realm of grounding natural language instructions for robotic systems, recent efforts have focused on mapping human commands to actionable representations that facilitate efficient task planning and execution. [2] interprets and executes natural language commands at varying levels of abstraction within a hierarchical planning framework to improve

grounding accuracy and [9] uses sentences to learn compositional reward function using lambda calculus. It's difficult to scale to complex instructions and environments because these approaches use predefined objects, spatial relations and language based features. [3] uses an adversarial learning framework where a discriminator differentiates between a predefined set of good (instruction, state) pairs and those generated by the current policy. The output of the discriminator is being used as a reward function to jointly optimize and improve the policy. The main difference between this approach and our approach is, it jointly trains the RL algorithms with linguistic features to improve policy, whereas in our method language is not a part of RL algorithm, rather (language, policy) pairs are being used to initialize policy for similar task.

Recent works [5, 6, 4] have introduced intermediate rewards as a means to guide reinforcement learning agents more effectively. However training these algorithm requires significant human intervention for generating expert trajectories. The main difference in methods proposed above and our method is using the language to directly train policies rather than explicitly learning reward function. Our work is closest to [7], where a set of base control policies along with task description in natural language can be used for task adaptation. However, [7] uses binary classification to sample task description for task adaption and our work makes use of CLIP [8] which is proven effective in contrastive learning between two different modalities and thus helps in selecting similar task description combined with policies.

# 3    Problem Description

Our objective is to efficiently initialize the policy network for a target task, guided by a natural language description, to accelerate subsequent training. To achieve this, we leverage the language instructions and policy networks from previously trained tasks.

We operate in a Markov Decision Process (MDP) setting where an agent sees the environment via set of states $s \in S$, acts on it with actions $a \in A$, observes the next states with transition probabilities $s_{t+1} \sim P(s_{t+1}|s_t, a_t)$, and receives a reward $r_t \sim R$ with $\gamma \in (0, 1]$ as discount factor. In addition to the standard MDP, our agent is provided with a natural language description of the tasks. Let $z$ represent the complete set of tasks to be learned within the environment.

---

**Algorithm 1**

---

$\alpha$ : sample $l$ descriptions from a set of z task descriptions
$\beta$ : sample $l'$ descriptions from a set of z task descriptions
$\tau(l)$ : takes sentence as input and gives corresponding vector embedding as output
$nn(\pi)$: neural network weights corresponding to the policy $\pi$
Step: 1 (Train Base Policies)
**for all** instruction $i \in \alpha$ **do**
    Train base policy $\pi_i$ until convergence
**end for**
Step: 2 (Identify similar language descriptions for transfer)
**for all** instruction $j \in \beta$ **do**
    **for all** instruction $i \in \alpha$ **do**
        $d_j^i \leftarrow \cos(\tau(j), \tau(i)) = \frac{\tau(j) \cdot \tau(i)}{\|\tau(j)\| \|\tau(i)\|}$
    **end for**
    $nn(\pi') \leftarrow \frac{\sum_i d_j^i nn(\pi_i)}{\sum_i d_j^i}$
**end for**
Initialize the policy network of target task as $nn(\pi')$.

---

# 4    Proposed Algorithms

We begin by discussing the algorithm that employs transfer solely on language-based similarity. Subsequently, we introduce our novel CLIP-inspired transfer technique, designed to address the limitations of the baseline approach.

## 4.1   Task adaptation using language similarity

The fundamental approach of this algorithm is to find similar language semantics between language description of given target task and each source (pre-trained) task. To accomplish this, we use cosine similarity as measure of similarity. Later, we use this similarity to initialize the policy parameters of target task by weighted average of source task policy parameters, weighted by normalized similarity scores. The pseudo-code of discussed algorithm is given in Algorithm 1.

Algorithm 1 captures the similar language semantics based on cosine similarity, but doesn't capture its association with policy. As a result, it does not capture the structure of the environment. In the next section, we describe a solution to mitigate this problem.

## 4.2 Task adaptation using CLIP

The central idea is that the natural language instruction and the optimal policy for a given task represent the same underlying concept, where the policy serves as the actionable realization of the task. Therefore, the embeddings of the language description and the policy representation should be aligned and closer in the shared representation space. Initially, we create the dataset $N$ $(l_i, \pi_i)$ by pairing instructions with corresponding trained optimal base policies. Subsequently, we train CLIP [8] on $N \times N$ $(l_i \times \pi_j)_{i,j=1}^N$ pairs. CLIP learns a multi-modal shared embedding space by training the policy encoder and instruction encoder, maximizing cosine similarity of policy and instruction embedding of the N diagonal pairs, simultaneously minimizing the cosine similarity of $N^2 - N$ non-diagonal pairs. Given the instruction for target task, we project the instruction embedding into shared embedding space and check the CLIP similarity with each instruction in pre-trained dataset. Subsequently, corresponding policies are weighted by normalized CLIP similarity score. We initialize the policy parameters of a given target task by summation of these weighted polices. The pseudo-code of our approach is given in Algorithm 2.

---

**Algorithm 2**

---

$\alpha$ : sample $l$ descriptions from a set of z descriptions
$\beta$ : sample $l'$ descriptions from a set of z descriptions
$\tau(l)$ : takes sentence as input and gives vector embedding representation as output
$nn(\pi)$: neural network weights corresponding to the policy $\pi$
$Proj(\tau(l))$, $Proj(nn(\pi))$: Projection of the language embeddings and policy weights, respectively, randomly initialized
Step: 1 (Train $N$ base policies)
**for all** instruction $i \in \alpha$ **do**
    Train base policy $\pi_i$ until convergence
**end for**
Step 2: (Prepare dataset for CLIP training and calculate loss)
    a. Compute the similarity matrix

$$S_{ij} = \langle Proj(\tau(l_i)), \ Proj(nn(\pi_j)) \rangle \quad \forall i, j \in \{1, 2, \ldots, N\} \text{ (where } \langle . \rangle \text{ is dot product)}$$

    b. Train Projections on the following loss function ($\delta$ is the temperature parameter)

$$L = -\frac{1}{2N} \sum_{i=1}^N \log \frac{\exp(S_{ii}/\delta)}{\sum_{j=1}^N \exp(S_{ij}/\delta)} - \frac{1}{2N} \sum_{j=1}^N \log \frac{\exp(S_{jj}/\delta)}{\sum_{i=1}^N \exp(S_{ij}/\delta)}$$

Step: 3 (Identify similar language descriptions for transfer)
**for all** instruction $j \in \beta$ **do**
    **for all** instruction $i \in \alpha$ **do**
        $d_j^i \leftarrow \cos(Proj(\tau(j)), Proj(\tau(i)) = \frac{Proj(\tau(j)) \cdot Proj(\tau(i)}{\|Proj(\tau(j))\|\|Proj(\tau(i))\|}$
    **end for**
    $nn(\pi') \leftarrow \frac{\sum_i d_j^i nn(\pi_i)}{\sum_i d_j^i}$
**end for**
Initialize the policy network of target task as $nn(\pi')$.

---

The main focus of Algorithm 2 is to bring both policy and it's corresponding natural language instruction into similar embedding space, so that the agent not only follows the policy which are closer to similar tasks but also understands the language semantics describing it. Algorithm 2 not only performs above par the Algorithm 1 in terms of computation time but also captures the policies of source tasks which are in the neighbourhood of given target task.

## 5 Experiments and Results

We consider grid environments of varying sizes: $8 \times 8$, $10 \times 10$, $15 \times 15$, and $25 \times 25$ for our experimental evaluation. In each environment, the agent is placed at a random position on the grid. The task for the agent is to reach the Goal in the fewest possible steps. For all grid configurations and both algorithms, we use four natural language instructions as base tasks: "top left first", "top left second", "top right first", and "top right second". Each base task is trained optimally. The
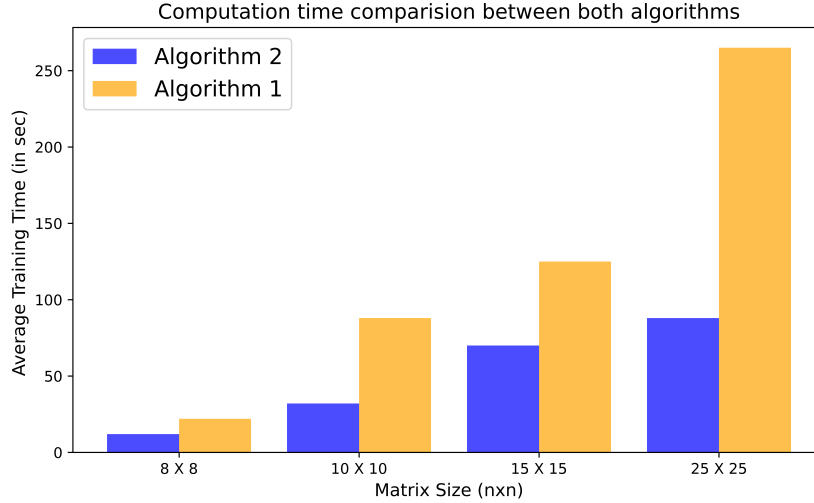
Figure 2: Transfer Performance comparison on grid world of varying sizes. We can observe that CLIP inspired Algorithm 2 outperforms Algorithm 1, which is simple language-based transfer.

target task for transfer is "top left third". These instructions represent specific goal positions in the grid. For instance, in a $10 \times 10$ grid, the instruction "top right first" places the goal at $(0, 9)$, while "top right second" places it at $(0, 8)$. For the target task, the goal is positioned at $(0, 7)$.

We run both algorithms for $10$ independent trials and report the average transfer performance in Figure 2. It can be observed that Algorithm 2 (inspired by CLIP) selects better policies, providing more effective initializations and thus enabling faster training of the target task compared to Algorithm 1. Moreover, as the size of the grid increases, the performance improves exponentially, demonstrating the scalability of the Algorithm 2.

## References

[1] Brenna D. Argall, Sonia Chernova, Manuela Veloso, and Brett Browning. A survey of robot learning from demonstration. *Robotics and Autonomous Systems*, 57(5):469–483, 2009.

[2] Dilip Arumugam, Siddharth Karamcheti, Nakul Gopalan, Lawson Wong, and Stefanie Tellex. Accurately and efficiently interpreting human-robot instructions of varying granularities. In *Robotics: Science and Systems XIII*, RSS2017. Robotics: Science and Systems Foundation, July 2017.

[3] Dzmitry Bahdanau, Felix Hill, Jan Leike, Edward Hughes, Arian Hosseini, Pushmeet Kohli, and Edward Grefenstette. Learning to understand goal specifications by modelling reward, 2019.

[4] Kate Baumli, Satinder Baveja, Feryal Behbahani, Harris Chan, Gheorghe Comanici, Sebastian Flennerhag, Maxime Gazeau, Kristian Holsheimer, Dan Horgan, Michael Laskin, Clare Lyle, Hussain Masoom, Kay McKinney, Volodymyr Mnih, Alexander Neitz, Dmitry Nikulin, Fabio Pardo, Jack Parker-Holder, John Quan, Tim Rocktäschel, Himanshu Sahni, Tom Schaul, Yannick Schroecker, Stephen Spencer, Richie Steigerwald, Luyu Wang, and Lei Zhang. Vision-language models as a source of rewards, 2024.

[5] Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale knowledge. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022.

[6] Prasoon Goyal, Scott Niekum, and Raymond J. Mooney. Using natural language for reward shaping in reinforcement learning, 2019.

[7] Matthias Hutsebaut-Buysse, Kevin Mets, and Steven Latré. Fast task-adaptation for tasks labeled using natural language in reinforcement learning, 2019.

[8] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021.

[9] Edward C. Williams, Nakul Gopalan, Mine Rhee, and Stefanie Tellex. Learning to parse natural language to grounded reward functions with weak supervision. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4430–4436, 2018.