

# STRUCTURED PREDICTIVE REPRESENTATIONS IN REINFORCEMENT LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Reinforcement Learning (RL) remains brittle in complex environments characterized by sparse rewards, partial observability, and subtask dependencies. Predictive state abstractions capture the environment’s underlying temporal structure and are crucial to overcoming these challenges. Yet, such methods often only focus on global one-step transitions and overlook local relationships between trajectories. This paper explores how capturing such relationships can enhance representation learning methods in RL. Our primary contribution is to show that incorporating a **Graph-Neural Network (GNN)** into the observation-predictive learning process improves sample efficiency and robustness to changes in size and distractors. Through experiments on the MiniGrid suite, we demonstrate that our GNN-based approach outperforms typical models that use **Multi-layer Perceptrons (MLPs)** in **sparse reward and partially-observable** environments where task decompositions are critical. These results highlight the value of structural inductive biases for generalization and adaptability, revealing how such mechanisms can bolster performance in RL.

## 1 INTRODUCTION

Environments with partial observability, sparse rewards, and dynamic changes frequently challenge Deep Reinforcement Learning (RL) algorithms, often rendering them brittle and sample-inefficient (Wang et al., 2019; Meng & Khushi, 2019; Lu et al., 2020; Tomar et al., 2023; Benjamins et al., 2023). Traditional RL methods struggle particularly in such complex environments due to the challenges of capturing long-term dependencies and relational structures between states. Learning representations of the state relevant to control offers a promising avenue to scale RL to complex scenarios. *State abstractions* in Markov Decision Processes (MDPs) (Dayan, 1993; Dean & Givan, 1997; Li et al., 2006) and *history abstractions* in Partially **Observable** MDPs (POMDPs) (Littman et al., 2001; Castro et al., 2009) improve data efficiency and generalization (Killian et al., 2017; Zhang et al., 2021). Consequently, numerous RL representation learning techniques have emerged in the last years (Castro et al., 2021; Schwarzer et al., 2021; Hansen-Estruch et al., 2022; Lan & Agarwal, 2023; Guo et al., 2020; Grill et al., 2020) making it an active area of research in RL.

*Self-prediction* has positioned itself as a prominent technique for learning state abstractions. It is a self-supervised mechanism that uses a latent model to predict the next latent state using the current latent state and action as inputs (Guo et al., 2019; 2020; Grill et al., 2020; Schwarzer et al., 2021; Lee et al., 2021). In doing so, it approximates the one-step transition structure in the latent space (Tang et al., 2023; Voelcker et al., 2024; Khetarpal et al., 2024). This objective is also connected to the objective to predict subsequent observations in POMDPs (Ni et al., 2024), allowing the agent to approximate the actual transition dynamics in the belief space (Schrittwieser et al., 2020; Subramanian et al., 2022). Real-world environments, however, often come with rich local structure as well (Mohan et al., 2024), which is usually overlooked by these methods.

This paper investigates how leveraging Graph Neural Networks (GNNs) (Battaglia et al., 2018) within a self-predictive framework can enhance representation learning in RL in **sparse reward and partially observable settings**. Specifically, we propose a method that

captures relationships between a batch of latent states generated by a history encoder. This approach enables the model to encode temporal and relational dependencies in the observation-prediction mechanism, improving the sample’s learning efficiency and robustness to environmental changes. In contrast to commonly used Multi-Layer Perceptron (MLP)-based methods, which often struggle with long-term dependencies and partial observability, GNNs excel at capturing relational structure between the latent states produced over time (see Figure 1).

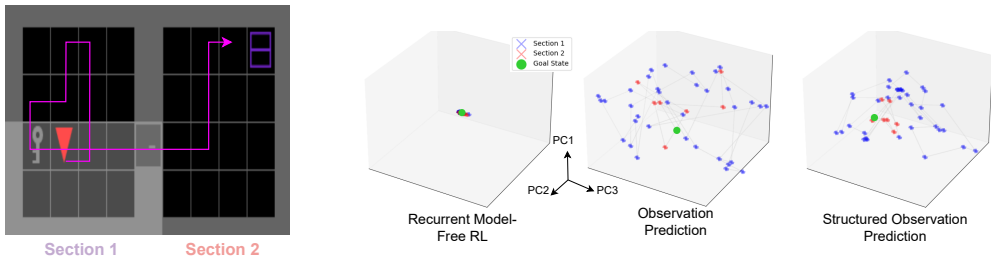


Figure 1: **Latent Space Representation.** A goal-reaching trajectory in MiniGrid-UnlockPickup-v0 mapped to a 3D PCA representation of the latent states generated by various belief encoders. States belonging to the first section are indicated in blue, while those in the second one are shown in red, with the goal state highlighted in green. Structured Observation Prediction captures the closeness of high-reward states (red) near the goal. In contrast, normal Observation Prediction reveals a less organized representation, indicating potential inefficiencies in recognizing rewarding states in this environment. This emphasizes the advantage of graph-based approaches for improved decision-making and performance in reinforcement learning tasks.

This paper’s main **contribution** is the introduction of a GNN-based observation-predictive model designed to operate on latent states generated by a history encoder. Unlike prior work that primarily focuses on spatial relationships (e.g., object-centric representations), our method targets temporal and relational dependencies in POMDPs. By relationally reasoning over trajectories, our method generalizes across variations in tasks. We validate our approach through experiments on a subset of navigation tasks in MiniGrid (Chevalier-Boisvert et al., 2023) that are particularly challenging for end-to-end observation prediction. Additionally, we demonstrate the robustness of our relational model in continually changing settings, showcasing its adaptability to distractors and environment size. Our results indicate that the GNN-based latent model outperforms MLP-based baselines, achieving superior performance in sparse-reward tasks and demonstrating better generalization to environmental variations.

## 2 BACKGROUND

In this section, we provide the necessary background to understand our approach. We briefly recap the fundamentals of RL, MDPs, and POMDPs, then delve deeper into state abstractions. Subsequently, we formally introduce self-predictive and Observation-Predictive (OP) abstractions, which we use to build our method.

### 2.1 MDPs, POMDPs AND REINFORCEMENT LEARNING

A discounted MDP (Puterman, 2014) is represented by a tuple  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, R, \gamma, \mu)$ . At each time step  $t$ , an agent observes the state  $s_t \sim \mathcal{S}$  of the environment and chooses an action  $a_t \sim \mathcal{A}$  using a policy  $\pi(a_t | s_t)$  to transition into a new state  $s_{t+1}$ . The dynamics govern the transitions function  $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ , and for each transition, the agent receives a reward according to the reward function  $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ . The agent’s objective is to maximize the expected cumulative discounted reward over an infinite horizon:

$$\max_{\pi} \mathbb{E}_{s_{t+1} \sim P(\cdot | s_t, a_t), a_t \sim \pi(\cdot | s_t)} \left[ \sum_{t=0}^{\infty} \gamma^{t-1} r_t \right] \quad (1)$$

where  $\gamma \in [0, 1]$  is the discount factor, and the starting state  $s_0$  is sampled from the initial state distribution  $s_0 \sim \mu(s_0)$ .

**Value-based methods** learn an optimal state-action value function  $Q^*(s, a)$ , the expected return after starting in state  $s$  and taking action  $a$ , by repeatedly performing two steps till convergence: (i) **Policy Evaluation**: computing a value function  $Q^\pi(s, a)$  quantifying the expected return after taking action  $a$  in state  $s$ :  $Q^\pi(s, a) = \mathbb{E}_\pi \left[ \sum_{i=t}^{\infty} \gamma^{i-t} r_{i+1} \mid s_t = s, a_t = a \right]$ ; and (ii) **Policy Improvement**: learning a new value function from which actions can be greedily selected to maximize  $Q^\pi(s, a)$ :  $\pi'(s_t) \in \arg \max_{a_t \in \mathcal{A}} Q^\pi(s_t, a_t)$

In many real-world scenarios, the agent cannot fully observe the environment. Such problems are modeled by POMDPs, defined as a tuple  $\mathcal{M}_O = (\mathcal{S}, \mathcal{O}, \mathcal{A}, P, R, \gamma, \mu)$ , where the agent has access to observations  $o \in \mathcal{O}$  based on the state  $s \in \mathcal{S}$ . It can utilize a history  $h_t := \{o_1, a_1, o_2, a_2, \dots, o_t\} \in \mathcal{H}$ , by concatenating observations and actions, where  $\mathcal{H}$  represents the set of all possible histories.

Since the agent lacks full observability, maintaining a belief state — a probability distribution over possible states given the history — is essential for optimal decision-making (Kaelbling et al., 1998). However, computing and updating such beliefs for high dimensional environments can quickly become intractable (Subramanian et al., 2022). Therefore, the agent requires a history encoder that maps the history to a Markovian representation  $\phi_O : \mathcal{H}_t \rightarrow \mathcal{Z}$ .

## 2.2 STATE ABSTRACTIONS, SELF-PREDICTION AND OBSERVATION-PREDICTION

A Q-function itself can be decomposed into two parts: (i) An encoder that  $\phi_{Q^*} : \mathcal{S} \rightarrow \mathcal{Z}$ , that maps the states to abstract states  $z \in \mathcal{Z}$ , also known as state abstractions (Li et al., 2006), or latent states (Gelada et al., 2019). (ii) A critic  $C : \mathcal{Z} \rightarrow \mathcal{Q}$  that predicts the Q-values using the latent state  $\mathcal{Z}$ . This decomposition requires the latent state-space  $\mathcal{Z}$  to have sufficient information to accurately predict  $Q^*$ , i.e. if  $\phi(s_i) = \phi(s_j)$ , then it must hold that  $Q^*(s_i) = Q^*(s_j)$ . We can additionally incentivize the latent state to predict the one-step transition probabilities (Equation (ZP)) and rewards (Equation (RP)), thereby preserving the environment’s dynamics in the latent space. Equation (ZP) ensures that the latent state is predictive of the subsequent latent state by mapping the joint latent state-action space to a distribution over the latent space  $\Delta(\mathcal{Z})$ . Consequently, such abstractions are **self-predictive abstractions**, learned using a latent model trained to predict the next latent state (Grill et al., 2020; Guo et al., 2020).

$$\exists P_z : \mathcal{Z} \times \mathcal{A} \rightarrow \Delta(\mathcal{Z}) \text{ s.t. } P(z_{t+1} \mid s_t, a_t) = P_z(z_{t+1} \mid \phi_L(s_t), a_t) \quad (\text{ZP})$$

$$\exists P_z : \mathcal{Z} \times \mathcal{A} \rightarrow \mathbb{R} \text{ s.t. } \mathbb{E}(r_{t+1} \mid h_t, a_t) = R_z(\phi_L(h_t), a_t) \quad (\text{RP})$$

For POMDPs, we can extend the state encoder to belief encoder  $\phi_O$  to produce a *history abstraction*  $z = \phi_O(h) \in \mathcal{Z}$ . This encoder satisfies as additional recurrent condition to ensure belief reconstruction:

$$\exists \psi_z : \mathcal{Z} \times \mathcal{A} \times \mathcal{O} \rightarrow \mathcal{Z} \text{ s.t. } \phi(h_{t+1}) = \psi_z(\phi_O(h_t), a_t, o_{t+1}) \quad (\text{Rec})$$

Furthermore, such abstractions should additionally satisfy a variant of Equation (ZP), called *Observation-prediction*, ensuring that the latent state along with the action is sufficient to predict the distribution over the subsequent observations (Equation (OP)):

$$\exists P_o : \mathcal{Z} \times \mathcal{A} \rightarrow \Delta(\mathcal{O}) \text{ s.t. } P(o_{t+1} \mid h_t, a_t) = P_o(o_{t+1} \mid \phi_O(h_t), a_t) \quad (\text{OP})$$

### 3 METHOD

In this section, we motivate and outline our method. We present the general idea of incorporating additional structure across batches of observations and the inter-trajectory transfer it enables. We then argue how capturing structure across batches is particularly beneficial for tasks with subtask decompositions, especially in a Sparse Reward environment. We then outline our architecture that incentivizes the belief encoder to produce such histories.

#### 3.1 RELATIONAL TASK DECOMPOSITION

Complex RL tasks often involve multiple subtasks. In sparse-reward MDPs, these subtasks are crucial but unrewarded steps, making learning challenging due to the delayed feedback. A vital requirement for credit assignment is to model the relationships across these subtasks to assign credit to the crucial state-action pairs. A state abstraction that preserves the optimal Q-value must enable the agent to disentangle latent states corresponding to these crucial ground states.

The intuition behind our approach is that trajectories corresponding to a single subtask exhibit correlations. In addition to the global one-step transition dynamics captured by self- and observation-predictive objectives, local structure among subtasks can be leveraged in the latent space. For example, consider the MDP shown in Section 3.1 where the agent must follow a goal-directed reward to the goal-state  $S_5$ . The reward includes a small cost per step to the agent and a large reward for reaching the goal. Therefore, the agent must discover the shortest path to reach the goal.

We highlight two example trajectories  $\tau_1 = \{S_1, R, S_3, U, S_5\}$  (illustrated in purple) and  $\tau_2 = \{S_2, R, S_4, U, S_5\}$  (shown in red). On reaching the goal  $S_5$ , it gets a reward of  $1 - k \times \text{n\_steps}$ , where  $k \in [0, 1)$ . It incentivizes the agent to reach the shortest path to the goal. The two trajectories involve two steps to the goal and accumulate the same return since they both comprise 2-steps to the goal. Although trained solely on data from  $\tau_1$ , a predictive model capable of capturing relational similarities between these trajectories can generalize to  $\tau_2$  by capturing local similarities between these trajectories. For instance, the relationship between  $S_3$  and  $S_5$  in  $\tau_1$  parallels the relationship between  $S_4$  and  $S_5$  in  $\tau_2$ .

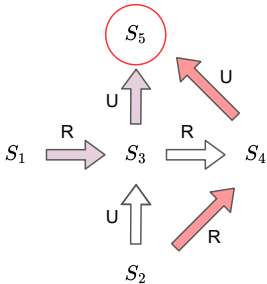


Figure 2: **Example MDP.** The agent must navigate to the goal  $S_5$  by maximizing a goal-conditioned reward and minimizing the cost per step. At the start of the episode, the agent can spawn in any of the other states  $\{S_1, S_2, S_3, S_4\}$ . From each state, it can either go right  $R$  or up  $U$ .

Let us extend this to the POMDP setting, where the agent does not directly observe the states. Instead, it receives partial observations corresponding to these states. The trajectories in this POMDP now correspond to histories of observations, actions, and rewards  $h_1 = \{o_1, a_1, r_1, \dots, o_5\}$  and  $h_2 = \{o_2, a_2, r_2, \dots, o_5\}$ . Here, the observations  $o_1, o_2, \dots$  are partial representations of the states  $S_1, S_2, \dots$ , and the goal is to navigate towards the final observation corresponding to  $S_5$ . Since the agent only observes part of the state, it must infer relationships and similarities between different observation sequences. As in the MDP case, the agent benefits from recognizing relational similarities between these histories to generalize across subtasks.

**Proposition 3.1.** *Let  $h_1, \dots, h_n \in \mathcal{H}$  be histories from similar subtasks in a POMDP, with corresponding next observations  $o'_1, \dots, o'_n \in \mathcal{O}$ . Let  $\phi : \mathcal{H} \rightarrow \mathcal{Z}$  be a Lipschitz continuous function with constant  $L_\phi > 0$ , mapping histories to embeddings  $z_i = \phi(h_i)$ . Let  $f : \mathcal{Z}^n \rightarrow \mathcal{O}$  be a Lipschitz continuous model with constant  $L_f > 0$ , predicting  $o'_{pred} = f(z_1, \dots, z_n)$ .*



Assume the histories  $h_i$  are similar, i.e.,  $d_{\mathcal{H}}(h_i, h_j) \leq \delta$  for all  $i, j$ , where  $d_{\mathcal{H}}$  measures the distance between histories.

Then, minimizing the squared error loss

$$\mathcal{L} = \|o'_{pred} - o'_i\|^2,$$

for any  $i$ , ensures that the prediction error is bounded:

$$\mathcal{L} \leq (L_f L_\phi n \delta + \epsilon_i)^2,$$

where  $\epsilon_i$  accounts for model approximation errors or inherent noise.

We sketch this proposition more intuitively by considering the trajectories in Figure 2 as histories. Since transitions from state  $S_3 \rightarrow S_5$  and  $S_4 \rightarrow S_5$  share a similar relational structure, the embeddings  $z_1 = \phi(\{S_3, U, S_5\})$  and  $z_2 = \phi(\{S_4, U, S_5\})$  will be close in the latent space. Training a model to minimize the loss  $\mathcal{L}$  by reasoning over both these trajectories ensures that the model generalizes between these subtasks, capturing the similarities between these histories. We do this using a GNN. Please refer to Appendix A.1 for a more detailed proof sketch.

### 3.2 OBSERVATION PREDICTION USING A GRAPH-BASED LATENT MODEL

Our method comprises three key components, illustrated in Figure 3:

1. **Encoder** ( $\phi$ ) that maps histories to latent representations  $z$ .
2. **Model** ( $\psi$ ) that captures relationships among history embeddings.
3. **RL network** ( $\pi_\theta$  or  $q_\theta$ ) that uses  $z$  for either learning a policy, or a  $Q$ - function depending on the method that we utilize.

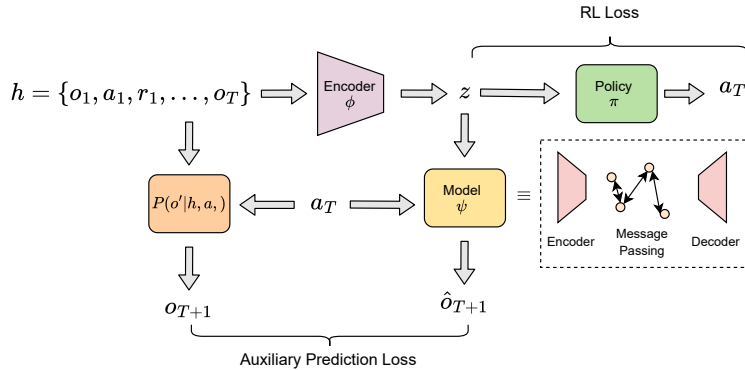


Figure 3: **Training Setup.** The LSTM generates embeddings using observation history, actions, and rewards, capturing temporal dependencies to create a belief state  $z$ . The policy network uses this to select the next action. For a value-based agent in a discrete action space, this would be a critic network that outputs values over discrete actions. Then, the algorithm greedily selects the action with the highest value. During optimization, the structured model - A GNN- reasons over a batch of latent states and corresponding actions to predict the subsequent observations. This is compared against the corresponding next observations to create the auxiliary prediction loss.

**Encoder and Policy Network.** The encoder  $\phi$  maps the history of observations to a latent state  $z = \phi(h_t)$ . In a POMDP, this is either a recurrent encoder (Subramanian et al., 2022) or possibly a transformer with a sufficiently large context window (Esslinger et al., 2022). Any RL agent can now use this latent state. In both these cases, the policy  $\pi(a_t|z_t)$  takes the latent state  $z_t$  as an input and outputs an action  $a_t$ . Value-based methods use a critic network that outputs values for each action for a given  $z$  and greedily selects the action with the maximum value.

**Graph construction.** The latent model  $\psi$  is a self-predictive model to enhance representation learning. To capture relational structure within the latent space, we consider a batch of latent states  $Z = [z_1, \dots, z_T]$  and corresponding actions  $\{a_1, \dots, a_T\}$ . We convert these actions to one-hot vectors and then concatenate them to form node features  $\{(z_1, a_1), \dots, (z_T, a_T)\}$ . Then, we construct a  $m$ -nearest neighbors graph on these with  $m = 4$  using the Euclidean distance between the node features.

**Message Passing.** After constructing the graph, the nodes with actions as attributes are passed through two message-passing layers. During this phase, each node in the graph updates its state by aggregating information from its neighboring nodes. Firstly, for each node, the features of its neighboring nodes are aggregated by concatenating the features of the source node  $x_i$  and the target node  $x_j$ . This concatenated vector is then passed through an MLP consisting of two fully connected layers with a ReLU activation function in between, transforming the combined features to capture more complex interactions. The result of this MLP is then used to update the target node’s features.

**Observation-Prediction and training.** After the message-passing steps, the updated node features are decoded to produce the final node representations. The output of the network has the same dimensionality as the flattened observation dimensions, and therefore, allows the graph to predict a batch of the subsequent observations  $\{\hat{o}_2, \dots, \hat{o}_{t+1}\}$  by reasoning across the batch of  $T$  observations and actions. The output of the GNN is then compared with the corresponding ground-truth observations  $\{o_2, \dots, o_{t+1}\}$  present in the buffer during training to create an auxiliary loss. This loss is jointly optimized along with the RL loss from the policy or critic network. As a result, we can train the encoder ( $\phi$ ), the model ( $\psi$ ), and the policy ( $\pi$ ) together during the optimization procedure.

$$\{\hat{o}_2, \dots, \hat{o}_{T+1}\} = \psi(\{[z_j, a_j]\}_{j=1}^T)$$

This output is trained using the Mean-Squared Error (MSE) loss between the predicted outputs  $\{\hat{o}'_1, \dots, \hat{o}_T\}$  and the actual next observation  $\{\hat{o}_1, \dots, \hat{o}_T\}$  sampled from the batch forming the representation learning auxiliary loss:

$$\mathcal{L}_{\text{OP}} = \sum_{t=1}^T \|\hat{o}_{t+1} - o_{t+1}\|^2$$

In principle, this objective is agnostic to the RL objective and, therefore, can be combined with any RL algorithm. We demonstrate an example of using our model with a policy-gradient algorithm in Algorithm 1.

**Reward Module.** For environments with multiple subtasks and sparse rewards, OP alone is insufficient (Ni et al., 2024). Instead, it must be combined with an explicit reward prediction using the latent state and action. For these environments, we utilize a two-layer MLP for such a module in addition to the latent model and train it using a phased training procedure, where the reward module is optimized separately from the end-to-end optimization of the bellman and representation learning loss. Instead, we interleave the optimization of the reward prediction from the representation learning modules by optimizing them one after the other.

## 4 EXPERIMENTS

In this section, we empirically investigate the effectiveness of our structured latent model. We employ the Minigrid suite (Chevalier-Boisvert et al., 2023), which consists of a series of mini-levels designed to test various aspects of learning and adaptation. The RL agent in our experiments is the R2D2 agent (Kapturowski et al., 2019), including a recurrent replay buffer with uniform sampling. Our hyperparameters can be found in A.2. In the following paragraphs, we divide our analysis based on specific research questions. Our presented results have been performed across 5 seeds with the aggregated IQMs (Agarwal et al., 2021).

**Algorithm 1** Training Procedure with a value-based agent

---

**Require:** Initialized encoder  $\phi$ , policy network  $\pi$ , auxiliary graph model  $\psi$

- 1: **while** not converged **do**
- 2:   **Collect Trajectories** using policy  $\pi(a_t | z_t)$ :
- 3:     Collect experiences  $\tau = \{(o_t, a_t, r_t, o_{t+1})\}$
- 4:     Compute  $z_t = \phi(o_t)$ ,  $z_{t+1} = \phi(o_{t+1})$
- 5:     Collect experiences  $\tau = \{(o_t, a_t, r_t, o_{t+1})\}$
- 6:     Compute  $z_t = \phi(o_t)$ ,  $z_{t+1} = \phi(o_{t+1})$
- 7:   **Compute RL Loss:**
- 8:     Compute target values:  $V_{target} = r_t + \gamma V(z_{t+1})$
- 9:     Estimate Q-values:  $Q(z_t, a_t) \leftarrow Q(z_t, a_t)$
- 10:     $\mathcal{L}_{RL} = \frac{1}{N} \sum_t (Q(z_t, a_t) - V_{target})^2$
- 11:   **Compute Observation-Prediction Loss:**
- 12:     Construct graphs  $G_t$  from  $z_t$
- 13:     Predict  $\hat{o}_{t+1} = \psi(G_t, a_t)$
- 14:      $\mathcal{L}_{OP} = \sum_t \|\hat{o}_{t+1} - o_{t+1}\|^2$
- 15:   **Update Parameters:**
- 16:      $\mathcal{L} = \mathcal{L}_{RL} + \lambda \mathcal{L}_{OP}$
- 17:     Minimize  $\mathcal{L}$  w.r.t.  $\phi$ ,  $\pi$ ,  $\psi$
- 18: **end while**

---

**Performance on static environments.** We first evaluate our model (Graph\_OP) on selected environments in Minigrid. Our baselines are the observation predictive algorithm (min-OP) and the observation and reward prediction algorithm (min-AIS) (Ni et al., 2024). min-OP follows the same pipeline but uses an MLP for the observation prediction task. The MLP predicts the subsequent observation for each latent state in a batch and does not use relational reasoning for the whole batch. min-AIS, on the other hand, extends min-OP by predicting the subsequent reward in addition to the observation, improving performance in environments where observation prediction alone is insufficient for effective representation learning. The critical distinction between our method and these baselines is how they process the latent observations and associated actions. In the MLP-based baselines, each combination of latent observation and action is processed independently to predict the subsequent observation. By contrast, our GNN-based approach first constructs a graph over all the latent observation-action pairs in the batch, applies message passing across the graph to model relational dependencies, and then predicts the subsequent observations for each element. Therefore, the performance difference between the baselines and our method primarily comes from this privileged reasoning. We consider environments with subtasks from the Minigrid suite challenging without representation learning and particularly challenging for observation prediction. Please note that R2D2, without representation learning, fails to accumulate notable returns in these environments, as indicated by the curves in Ni et al. (2024). Moreover, we run each environment until the baselines demonstrate convergent behavior. Based on the learning curves provided by Ni et al. (2024), we narrow down the environments to the following four static ones:

1. MiniGrid-DoorKey-8x8-v0: The agent must pick up a key to unlock a door and reach the green goal in a  $8 \times 8$  grid.
2. MiniGrid-ObstructedMaze-1D1-v0: A blue ball is hidden in a maze with two rooms. A locked door separates the two rooms, and a ball obstructs the doors. The keys are hidden in boxes.
3. MiniGrid-KeyCorridorS3R2-v0: The agent has to pick up an object behind a locked door. The key is hidden in another room, and the agent has to explore the environment to find it.
4. MiniGrid-UnlockPickup-v0: The agent must pick up a box behind a locked door in another room.

These environments share the commonality of subtasks the agent needs to solve before reaching the goal. Apart from the DoorKey environment, all others require additional reward

prediction due to the sparsity of the reward in the original task. Consequently, we incorporate an additional reward-prediction module with our graph prediction (**Graph\_AIS**).

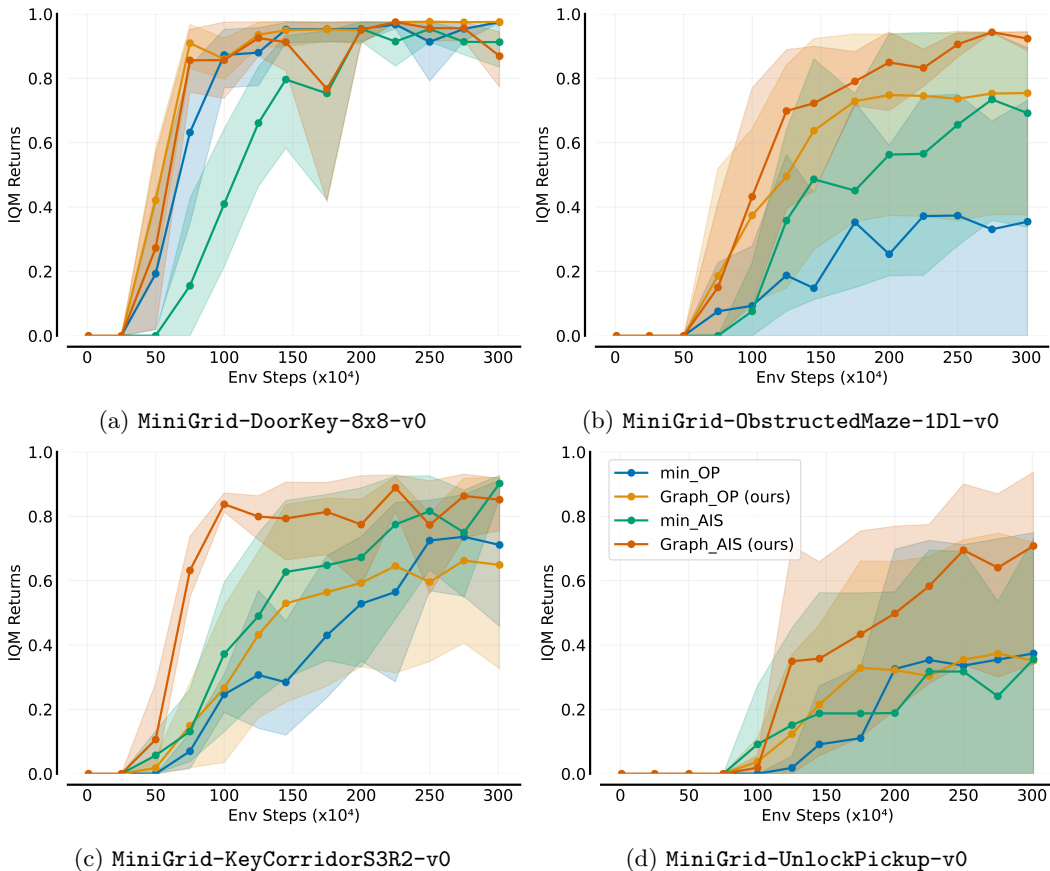


Figure 4: IQM and quartiles of Performances on static environments.

Our results are presented in Figure 4. Overall, the **Graph**-based representation learning methods outperform the MLP-based techniques in most cases. For environments where observation prediction struggles with long-term dependencies, the combination of **Graph**-based observation prediction and reward prediction – **Graph\_AIS** – consistently outperforms the baselines. This reiterates the inefficiencies of pure observation prediction in such environments since the reward is highly sparse in these subtasks.

**Adapting to environment changes.** A crucial outcome of Proposition 3.1 would be the ability of our method to extrapolate the learned prediction across environmental changes insofar as these changes share some similarity with data seen already. We investigate this by creating a scenario where an agent must continually adapt to environmental variations. We introduce changes to **MiniGrid-DoorKey-8x8-v0** by changing: (i) **Number of keys:** We introduce distractions in the form of additional colorless keys, forcing the agent to focus on the colored key. The number of distractors remains constant for each episode, but their location changes after the reset. (ii) **Size:** We periodically increase the size of the environment to investigate how well the agent adapts to the increase in the number of states.

Figure 5 shows the performance of **Graph\_OP** against **min\_OP** for different types of changes. Figure 5(a) demonstrates the agent’s performance when distractors are added after 800K steps, and Figure 5(b) shows the adaptation to increase in size after 1M steps. We introduce additional dimensions of hardness by combining these changes. Figure 5(c) shows the scenario in which the grid increases in size every 1M step, and a distractor is simultaneously added. Finally, Figure 5(d) shows the scenario in which the agents must adapt to a new distractor every 600K step and a size increment every 1M step in the bottom right figure.

As expected, both methods’ performance generally degrades when changes occur, and recovery from these changes becomes increasingly difficult as we increase hardness. As a result, in Figure 5(d), neither method has enough time to return to stable performance. In most of these scenarios, **Graph\_OP** remains consistently more robust performance and outperforms **min\_OP**. The impact of distractions seems more pronounced than size, as shown in Figure 5(a).

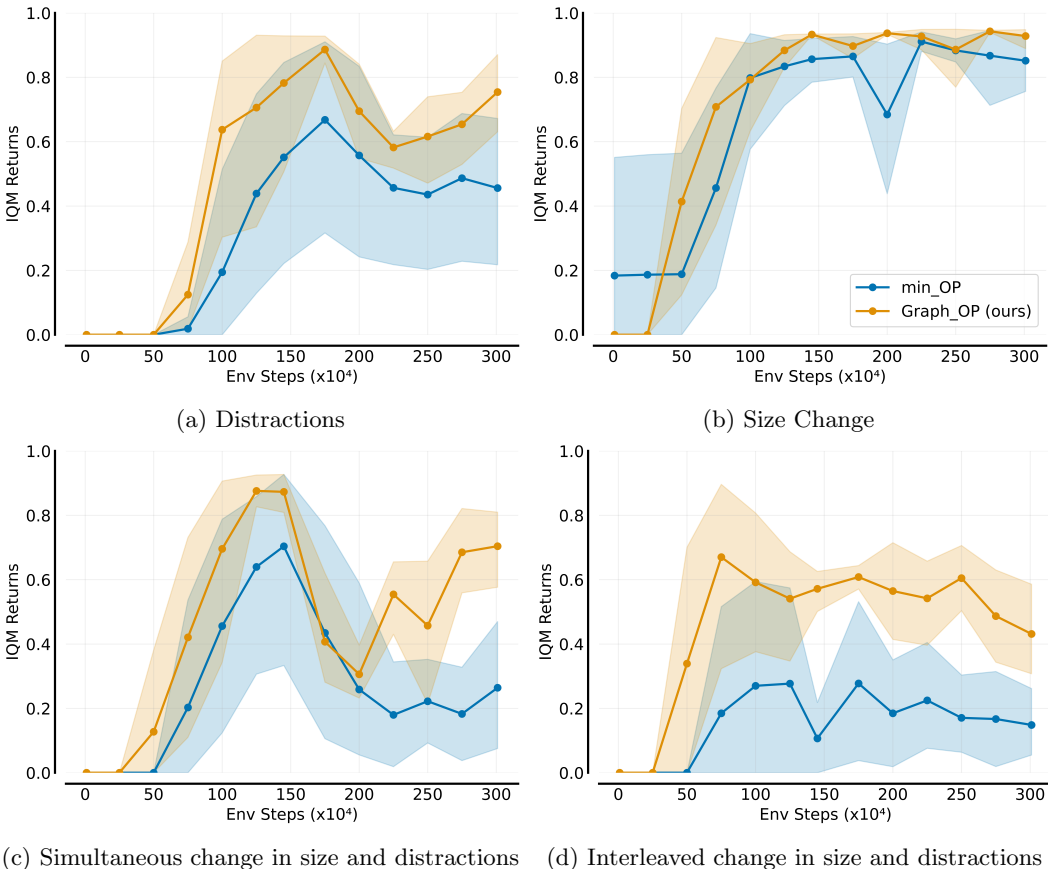


Figure 5: Performances on Dynamic Variations of MiniGrid-DoorKey-8x8-v0.

Compound changes particularly impact both methods since the size change forces the agent to explore more, while the distractors force the agent to focus on the right kind of key. Given that in DoorKey, the agent has to traverse a sub-goal of getting to a key before reaching a door and then going to the goal, changing the size and adding distractors together degrades performance faster. In both cases, the graph-based agent **Graph\_OP** is more robust to the changes than the MLP baseline. This highlights the particular advantage relational inductive bias offers: it allows the state representations to model relationships between trajectories and the one-step temporal consistency of self-prediction.

## 5 RELATED WORK

Our work touches upon three crucial areas in RL: *Abstractions, GNNs in RL, and incorporating structure in RL*, summarized below.

**State and History Abstractions in RL.** State abstractions constitute an active area in RL, and a complete categorization of approaches is beyond the scope of this work. Model-irrelevance has been studied under a variety of techniques, such as bi-simulation (Ferns et al., 2004; Gelada et al., 2019; Castro et al., 2021; Hansen-Estruch et al., 2022; Lan &

Agarwal, 2023), variational inference (Eysenbach et al., 2021; Ghugare et al., 2023), and successor features (Dayan, 1993; Barreto et al., 2017; Borsa et al., 2019; Lehnert & Littman, 2020; Scarpellini et al., 2024). Self-predictive representations have been a separate line of work (Guo et al., 2020; Grill et al., 2020; Schrittwieser et al., 2020; Schwarzer et al., 2021; Hansen et al., 2022; Ghugare et al., 2023; Zhao et al., 2023) with increasing interest in understanding how these objectives behave (Tang et al., 2023; Ni et al., 2024; Fang & Stachenfeld, 2024; Voelcker et al., 2024; Khetarpal et al., 2024). Observation predictive representations have been used to formulate belief states (Kaelbling et al., 1998; Wayne et al., 2018; Hafner et al., 2019; Han et al., 2020; Lee et al., 2020) and predictive state representations (Littman et al., 2001; Zhang et al., 2019), and are also related to observation reconstruction objectives commonly used for improving sample efficiency Yarats et al. (2021). Our work adds to this line of work by exploring how the self-predictive objective can capture relational structure in the latent space.

**Structure in RL.** Structural decompositions can be useful as inductive biases for various purposes (Mohan et al., 2024). Our work assumes a relational decomposition in joint state-action space. Such assumptions have previously been applied through modeling frameworks such as Relational MDPs (Dzeroski et al., 2001; Guestrin et al., 2003) and object-oriented MDPs (Diuk et al., 2008). However, we neither model entities in the environment separately nor handcraft any form of first-order representation in the value function (Guestrin et al., 2003; Fern et al., 2006; Joshi & Kharon, 2011). Instead, we reason across trajectories using a GNN to model relationships.

**GNNs in RL.** GNNs have increasingly been used in RL, such as modeling environments (Chen et al., 2020; Chadalapaka et al., 2023), agent’s morphology in embodied control (Wang et al., 2018; Oliva et al., 2022), relationships between different action sets (Jain et al., 2021), and concurrent policy optimization method (Wang & van Hoof, 2022). We share similarities to methods that use GNNs as structured models, used for applications such as learning the latent transition dynamics in simple manipulation tasks (Kipf et al., 2020), the dynamics of joints of physical bodies (Sanchez-Gonzalez et al., 2020), obtaining object-centric representations from images and RRT planners (Driess et al., 2022), or computing intrinsic reward and online planning (Sancaktar et al., 2022). We add to this line of work by using GNNs for observation-prediction. Although Transformers have also been used for learning state representations (Zhu et al., 2022) and state-action representations (Zheng et al., 2024), they require substantial data and computational resources, often making them less practical in data-scarce RL settings. In contrast, GNNs effectively leverage structural properties in relational tasks, providing an efficient alternative for relational reasoning in reinforcement learning.

## 6 CONCLUSION AND FUTURE WORK

Using a structured latent model to investigate the impact of relational inductive biases, we incorporated a GNN to capture the similarity between the latent space belief representations produced by a recurrent encoder. Our experiments on a relevant subset of Minigrid tasks demonstrated that agents utilizing this latent space exhibit improved performance and the learned representations tend to be more robust to changes in size and against added distractions. Although effective, our approach has been evaluated only on discrete action spaces and requires further investigations on continuous action spaces in environments such as robotic control (Freeman et al., 2021; Todorov et al., 2012), and on more complicated navigation topologies such as those found in Cobbe et al. (2020); Samvelyan et al. (2021). Additionally, we want to incorporate more algorithms since the current framework is agnostic to the RL algorithm. Finally, we want to extend our method to 3D point clouds to capture granular structure. Despite these limitations, our current findings offer a foundation for future research, and addressing these challenges will be crucial to advancing the capabilities of graph-based latent models in RL.



## REFERENCES

- 540  
541  
542 R. Agarwal, M. Schwarzer, P. Samuel Castro, A. C. Courville, and M. G. Bellemare.  
543 Deep reinforcement learning at the edge of the statistical precipice. In M. Ranzato,  
544 A. Beygelzimer, K. Nguyen, P. Liang, J. Vaughan, and Y. Dauphin (eds.), *Proceedings of*  
545 *the 35th International Conference on Advances in Neural Information Processing Systems*  
546 *(NeurIPS'21)*. Curran Associates, 2021.
- 547 A. Barreto, W. Dabney, R. Munos, J. Hunt, T. Schaul, D. Silver, and H. Hasselt. Successor  
548 features for transfer in reinforcement learning. In I. Guyon, U. von Luxburg, S. Bengio,  
549 H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Proceedings of the*  
550 *31st International Conference on Advances in Neural Information Processing Systems*  
551 *(NeurIPS'17)*. Curran Associates, 2017.
- 552 P. Battaglia, J. Hamrick, V. Bapst, A. Sanchez-Gonzalez, V. Zambaldi, M. Malinowski,  
553 A. Tacchetti, D. Raposo, A. Santoro, R. Faulkner, C. Gülçehre, H. Song, A. Ballard,  
554 J. Gilmer, G. Dahl, A. Vaswani, K. Allen, C. Nash, V. Langston, C. Dyer, N. Heess,  
555 D. Wierstra, P. Kohli, M. Botvinick, O. Vinyals, Y. Li, and R. Pascanu. Relational  
556 inductive biases, deep learning, and graph networks. *CoRR*, abs/1806.01261, 2018. URL  
557 <http://arxiv.org/abs/1806.01261>.
- 558 C. Benjamins, T. Eimer, F. Schubert, A. Mohan, S. Döhler, A. Biedenkapp, B. Rosenhan,  
559 F. Hutter, and M. Lindauer. Contextualize me – the case for context in reinforcement  
560 learning. *Transactions on Machine Learning Research*, 2023.
- 561  
562 D. Borsa, A. Barreto, J. Quan, D. Mankowitz, H. van Hasselt, R. Munos, D. Silver, and  
563 T. Schaul. Universal successor features approximators. In *Proceedings of the International*  
564 *Conference on Learning Representations (ICLR'19)*, 2019. Published online: [iclr.cc](https://iclr.cc).
- 565 P. Castro, P. Panangaden, and D. Precup. Equivalence relations in fully and partially  
566 observable markov decision processes. In *IJCAI 2009, Proceedings of the 21st International*  
567 *Joint Conference on Artificial Intelligence*, 2009.
- 568 P. Castro, T. Kastner, P. Panangaden, and M. Rowland. MICo: Improved representa-  
569 tions via sampling-based state similarity for markov decision processes. In M. Ranzato,  
570 A. Beygelzimer, K. Nguyen, P. Liang, J. Vaughan, and Y. Dauphin (eds.), *Proceedings of*  
571 *the 35th International Conference on Advances in Neural Information Processing Systems*  
572 *(NeurIPS'21)*. Curran Associates, 2021.
- 573  
574 V. Chadalapaka, V. Ustun, and L. Liu. Leveraging graph networks to model environments in  
575 reinforcement learning. In *Proceedings of the Thirty-Sixth International Florida Artificial*  
576 *Intelligence Research Society Conference (FLAIRS'23)*, 2023.
- 577 C. Chen, S. Hu, P. Nikdel, G. Mori, and M. Savva. Relational graph learning for crowd  
578 navigation. In *Proceedings of the 2020 IEEE/RSJ International Conference on Intelligent*  
579 *Robots and Systems (IROS'20)*, 2020.
- 580 M. Chevalier-Boisvert, B. Dai, M. Towers, R. de Lazcano, L. Willems, S. Lahlou, S. Pal,  
581 P. Castro, and J. Terry. Minigrid & miniworld: Modular & customizable reinforcement  
582 learning environments for goal-oriented tasks. *CoRR*, abs/2306.13831, 2023.
- 583  
584 K. Cobbe, C. Hesse, J. Hilton, and J. Schulman. Leveraging procedural generation to  
585 benchmark reinforcement learning. In H. Daume III and A. Singh (eds.), *Proceedings of the*  
586 *37th International Conference on Machine Learning (ICML'20)*, volume 98. Proceedings  
587 of Machine Learning Research, 2020.
- 588 P. Dayan. Improving generalization for temporal difference learning: The successor represen-  
589 tation. *Neural Comput.*, 5(4):613–624, 1993.
- 590  
591 T. Dean and R. Givan. Model minimization in markov decision processes. In *Proceedings*  
592 *of the Fourteenth National Conference on Artificial Intelligence and Ninth Innovative*  
593 *Applications of Artificial Intelligence Conference, AAAI 97, IAAI 97, July 27-31, 1997,*  
*Providence, Rhode Island, USA, 1997.*

- 594 C. Diuk, A. Cohen, and M. Littman. An object-oriented representation for efficient rein-  
595 forcement learning. In W. Cohen, A. McCallum, and S. Roweis (eds.), *Proceedings of the*  
596 *25th International Conference on Machine Learning (ICML'08)*. Omnipress, 2008.
- 597 D. Driess, Z. Huang, Y. Li, R. Tedrake, and M. Toussaint. Learning multi-object dynamics  
598 with compositional neural radiance fields. In *Conference on Robot Learning (CoRL'22)*,  
599 2022.
- 600 S. Dzeroski, L. Raedt, and K. Driessens. Relational reinforcement learning. *Machine Learning*,  
601 43(1/2):7–52, 2001.
- 602 K. Esslinger, R. Platt, and C. Amato. Deep transformer q-networks for partially observable  
603 reinforcement learning. *CoRR*, abs/2206.01078, 2022.
- 604 B. Eysenbach, R. Salakhutdinov, and S. Levine. Robust predictable control. In M. Ranzato,  
605 A. Beygelzimer, K. Nguyen, P. Liang, J. Vaughan, and Y. Dauphin (eds.), *Proceedings of*  
606 *the 35th International Conference on Advances in Neural Information Processing Systems*  
607 *(NeurIPS'21)*. Curran Associates, 2021.
- 608 C. Fang and K. Stachenfeld. Predictive auxiliary objectives in deep RL mimic learning in  
609 the brain. In *iclr24*, 2024.
- 610 A. Fern, S. Yoon, and R. Givan. Approximate policy iteration with a policy language bias:  
611 Solving relational markov decision processes. *Journal of Artificial Intelligence Research*,  
612 25:75–118, 2006.
- 613 N. Ferns, P. Panangaden, and D. Precup. Metrics for finite markov decision processes.  
614 In R. Holte and A. Howe (eds.), *Proceedings of the Nineteenth National Conference on*  
615 *Artificial Intelligence (AAAI'04)*. AAAI Press, 2004.
- 616 C. Freeman, E. Frey, A. Raichuk, S. Girgin, I. Mordatch, and O. Bachem. Brax - A  
617 differentiable physics engine for large scale rigid body simulation. In *Proceedings of the*  
618 *Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS*  
619 *Datasets and Benchmarks 2021*, 2021.
- 620 C. Gelada, S. Kumar, J. Buckman, O. Nachum, and M. Bellemare. DeepMDP: Learn-  
621 ing continuous latent space models for representation learning. In K. Chaudhuri and  
622 R. Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine*  
623 *Learning (ICML'19)*, volume 97. Proceedings of Machine Learning Research, 2019.
- 624 R: Ghugare, H. Bharadhwaj, B. Eysenbach, S. Levine, and R. Salakhutdinov. Simplifying  
625 model-based RL: learning representations, latent-space models, and policies with one  
626 objective. In *International Conference on Learning Representations (ICLR'23)*, 2023.  
627 Published online: [iclr.cc](https://arxiv.org/abs/2302.08801).
- 628 J. Grill, F. Strub, F. Altché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch, B. Pires,  
629 Z. Guo, M. Azar, B. Piot, K. Kavukcuoglu, R. Munos, and M. Valko. Bootstrap your  
630 own latent - A new approach to self-supervised learning. In H. Daume III and A. Singh  
631 (eds.), *Proceedings of the 37th International Conference on Machine Learning (ICML'20)*,  
632 volume 98. Proceedings of Machine Learning Research, 2020.
- 633 C. Guestrin, D. Koller, C. Gearhart, and N. Kanodia. Generalizing plans to new environments  
634 in relational MDPs. In G. Gottlob and T. Walsh (eds.), *Proceedings of the 18th International*  
635 *Joint Conference on Artificial Intelligence (IJCAI'03)*, 2003.
- 636 Y. Guo, J. Choi, M. Moczulski, S. Bengio, M. Norouzi, and H. Lee. Efficient exploration with  
637 self-imitation learning via trajectory-conditioned policy. *arXiv preprint arXiv:1907.10247*,  
638 2019.
- 639 Z. Guo, B. Pires, B. Piot, J. Grill, F. Altché, R. Munos, and M. Azar. Bootstrap latent-  
640 predictive representations for multitask reinforcement learning. In H. Daume III and  
641 A. Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*  
642 *(ICML'20)*, volume 98. Proceedings of Machine Learning Research, 2020.

- 648 D. Hafner, T. Lillicrap, I. Fischer, R. Villegas, D. Ha, H. Lee, and J. Davidson. Learn-  
649 ing latent dynamics for planning from pixels. In K. Chaudhuri and R. Salakhutdinov  
650 (eds.), *Proceedings of the 36th International Conference on Machine Learning (ICML'19)*,  
651 volume 97. Proceedings of Machine Learning Research, 2019.
- 652 D. Han, K. Doya, and J. Tani. Variational recurrent models for solving partially observable  
653 control tasks. In *iclr20*, 2020.
- 655 N. Hansen, H. Su, and X. Wang. Temporal difference learning for model predictive control. In  
656 K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvári, G. Niu, and S. Sabato (eds.), *Proceedings*  
657 *of the 39th International Conference on Machine Learning (ICML'22)*, volume 162 of  
658 *Proceedings of Machine Learning Research*. PMLR, 2022.
- 659 P. Hansen-Estruch, A. Zhang, A. Nair, P. Yin, and S. Levine. Bisimulation makes analog-  
660 ies in goal-conditioned reinforcement learning. In K. Chaudhuri, S. Jegelka, L. Song,  
661 C. Szepesvári, G. Niu, and S. Sabato (eds.), *Proceedings of the 39th International Confer-*  
662 *ence on Machine Learning (ICML'22)*, volume 162 of *Proceedings of Machine Learning*  
663 *Research*. PMLR, 2022.
- 665 A. Jain, N. Kosaka, K. Kim, and J. Lim. Know your action set: Learning action relations  
666 for reinforcement learning. In M. Meila and T. Zhang (eds.), *Proceedings of the 38th*  
667 *International Conference on Machine Learning (ICML'21)*, volume 139 of *Proceedings of*  
668 *Machine Learning Research*. PMLR, 2021.
- 669 S. Joshi and R. Khordon. Probabilistic relational planning with first order decision diagrams.  
670 *Journal of Artificial Intelligence Research*, 41:231–266, 2011.
- 672 L. Kaelbling, M. Littman, and A. Cassandra. Planning and acting in partially observable  
673 stochastic domains. *Artificial Intelligence*, 1998.
- 674 S. Kapturowski, G. Ostrovski, J. Quan, R. Munos, and W. Dabney. Recurrent experience  
675 replay in distributed reinforcement learning. In *Proceedings of the International Conference*  
676 *on Learning Representations (ICLR'19)*, 2019. Published online: [iclr.cc](https://arxiv.org/abs/1905.00981).
- 678 K. Khetarpal, Z. Guo, B. Pires, Y. Tang, C. Lyle, M. Rowland, N. Heess, D. Borsa, A. Guez,  
679 and W. Dabney. A unifying framework for action-conditional self-predictive reinforcement  
680 learning. *CoRR*, abs/2406.02035, 2024. doi: 10.48550/ARXIV.2406.02035.
- 681 T. Killian, S. Daulton, F. Doshi-Velez, and G. Konidaris. Robust and efficient transfer  
682 learning with hidden parameter markov decision processes. In I. Guyon, U. von Luxburg,  
683 S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Proceedings of*  
684 *the 31st International Conference on Advances in Neural Information Processing Systems*  
685 *(NeurIPS'17)*. Curran Associates, 2017.
- 686 T. Kipf, E. van der Pol, and M. Welling. Contrastive learning of structured world models. In  
687 *Proceedings of the 8th International Conference on Learning Representations (ICLR'20)*,  
688 2020.
- 690 C. Lan and R. Agarwal. Revisiting bisimulation: A sampling-based state similarity pseudo-  
691 metric. In *The First Tiny Papers Track at the 11th International Conference on Learning*  
692 *Representations (ICLR'23)*, 2023.
- 693 A. Lee, A. Nagabandi, P. Abbeel, and S. Levine. Stochastic latent actor-critic: Deep rein-  
694 forcement learning with a latent variable model. In H. Larochelle, M. Ranzato, R. Hadsell,  
695 M.-F. Balcan, and H. Lin (eds.), *Proceedings of the 34th International Conference on*  
696 *Advances in Neural Information Processing Systems (NeurIPS'20)*. Curran Associates,  
697 2020.
- 699 J. Lee, Q. Lei, N. Saunshi, and J. Zhuo. Predicting what you already know helps: Provable  
700 self-supervised learning. In M. Ranzato, A. Beygelzimer, K. Nguyen, P. Liang, J. Vaughan,  
701 and Y. Dauphin (eds.), *Proceedings of the 35th International Conference on Advances in*  
*Neural Information Processing Systems (NeurIPS'21)*. Curran Associates, 2021.

- 702 L. Lehnert and M. Littman. Successor features combine elements of model-free and model-  
703 based reinforcement learning. *J. Mach. Learn. Res.*, 2020.  
704
- 705 L. Li, T. Walsh, and M. Littman. Towards a unified theory of state abstraction for mdps.  
706 In *Proceedings of the International Symposium on Artificial Intelligence and Mathematics*  
707 (*AI&M'06*), 2006.
- 708 M. Littman, R. Sutton, and S. Singh. Predictive representations of state. In *Proceedings of*  
709 *the 15th International Conference on Advances in Neural Information Processing Systems*  
710 (*NeurIPS'01*), 2001.  
711
- 712 M. Lu, Z. Shahn, D. Sow, F. Doshi-Velez, and L. Lehman. Is deep reinforcement learning  
713 ready for practical applications in healthcare? a sensitivity analysis of duel-ddqn for  
714 hemodynamic management in sepsis patients. In *Proceedings of the American Medical*  
715 *Informatics Association Annual Symposium (AMIA '20)*, 2020.
- 716 T. Meng and M. Khushi. Reinforcement learning in financial markets. *Data*, 4(3):110, 2019.  
717
- 718 A. Mohan, A. Zhang, and M. Lindauer. Structure in deep reinforcement learning: A survey  
719 and open problems. *Journal of Artificial Intelligence Research*, 79, 2024.
- 720 T. Ni, B. Eysenbach, E. Seyedsalehi, M. Ma, C. Gehring, A. Mahajan, and P. Bacon. Bridging  
721 state and history representations: Understanding self-predictive rl. In *Proceedings of the*  
722 *12th International Conference on Learning Representations (ICLR'24)*, 2024.  
723
- 724 M. Oliva, S. Banik, J. Josifovski, and A. Knoll. Graph neural networks for relational inductive  
725 bias in vision-based deep reinforcement learning of robot control. In *International Joint*  
726 *Conference on Neural Networks, IJCNN 2022, Padua, Italy, July 18-23, 2022*, pp. 1–9,  
727 2022.
- 728 M. Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John  
729 Wiley & Sons, 2014.  
730
- 731 M. Samvelyan, R. Kirk, V. Kurin, J. Parker-Holder, M. Jiang, E. Hambro, F. Petroni,  
732 H. Kuttler, E. Grefenstette, and T. Rocktäschel. Minihack the planet: A sandbox for  
733 open-ended reinforcement learning research. In M. Ranzato, A. Beygelzimer, K. Nguyen,  
734 P. Liang, J. Vaughan, and Y. Dauphin (eds.), *Proceedings of the 35th International*  
735 *Conference on Advances in Neural Information Processing Systems (NeurIPS'21)*. Curran  
736 Associates, 2021.
- 737 C. Sancaktar, S. Blaes, and G. Martius. Curious exploration via structured world models  
738 yields zero-shot object manipulation. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave,  
739 K. Cho, and A. Oh (eds.), *Proceedings of the 36th International Conference on Advances*  
740 *in Neural Information Processing Systems (NeurIPS'22)*. Curran Associates, 2022.  
741
- 742 A. Sanchez-Gonzalez, J. Godwin, T. Pfaff, R. Ying, J. Leskovec, and P. Battaglia. Learning  
743 to simulate complex physics with graph networks. In H. Daume III and A. Singh  
744 (eds.), *Proceedings of the 37th International Conference on Machine Learning (ICML'20)*,  
745 volume 98. Proceedings of Machine Learning Research, 2020.
- 746 G. Scarpellini, K. Konyushkova, C. Fantacci, T. Le Paine, Y. Chen, and M. Denil.  $\pi 2vec$ :  
747 Policy representations with successor features. In *International Conference on Learning*  
748 *Representations (ICLR'24)*, 2024. Published online: [iclr.cc](https://iclr.cc).
- 749 J. Schrittwieser, I. Antonoglou, T. Hubert, K. Simonyan, L. Sifre, S. Schmitt, A. Guez,  
750 E. Lockhart, D. Hassabis, T. Graepel, T. Lillicrap, and D. Silver. Mastering atari, go,  
751 chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.  
752
- 753 M. Schwarzer, A. Anand, R. Goel, R. Hjelm, A. Courville, and P. Bachman. Data-efficient  
754 reinforcement learning with self-predictive representations. In *Proceedings of the Inter-*  
755 *national Conference on Learning Representations (ICLR'21)*, 2021. Published online:  
[iclr.cc](https://iclr.cc).

- 756 J. Subramanian, A. Sinha, R. Seraj, and A. Mahajan. Approximate information state for  
757 approximate planning and reinforcement learning in partially observed systems. *Journal*  
758 *of Machine Learning Research*, 2022.
- 759 Y. Tang, Z. Daniel Guo, P. Richemond, B. Pires, Y. Chandak, R. Munos, M. Rowland,  
760 M. Gheshlaghi Azar, C. Lan, C. Lyle, A. György, S. Thakoor, W. Dabney, B. Piot,  
761 D. Calandriello, and M. Valko. Understanding self-predictive learning for reinforcement  
762 learning. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett  
763 (eds.), *Proceedings of the 40th International Conference on Machine Learning (ICML’23)*,  
764 volume 202 of *Proceedings of Machine Learning Research*. PMLR, 2023.
- 765 E. Todorov, T. Erez, and Y. Tassa. MuJoCo: A physics engine for model-based control. In  
766 *International Conference on Intelligent Robots and Systems (IROS’12)*, pp. 5026–5033.  
767 ieeecis, IEEE, 2012.
- 768 M. Tomar, U. Mishra, A. Zhang, and M. Taylor. Learning representations for pixel-based  
769 control: What matters and why? *Trans. Mach. Learn. Res.*, 2023, 2023.
- 771 C. Voelcker, T. Kastner, I. Gilitschenski, and A. Farahmand. When does self-prediction  
772 help? understanding auxiliary tasks in reinforcement learning. *Reinforcement Learning*  
773 *Journal*, 2024.
- 774 Q. Wang and H. van Hoof. Model-based meta reinforcement learning using graph structured  
775 surrogate models and amortized policy search. In K. Chaudhuri, S. Jegelka, L. Song,  
776 C. Szepesvári, G. Niu, and S. Sabato (eds.), *Proceedings of the 39th International Confer-*  
777 *ence on Machine Learning (ICML’22)*, volume 162 of *Proceedings of Machine Learning*  
778 *Research*. PMLR, 2022.
- 779 T. Wang, R. Liao, J. Ba, and S. Fidler. Nervenet: Learning structured policy with graph  
780 neural networks. In *Proceedings of the 6th International Conference on Learning Repr-*  
781 *esentations (ICLR’18)*, 2018.
- 782 T. Wang, X. Bao, I. Clavera, J. Hoang, Y. Wen, E. Langlois, S. Zhang, G. Zhang, P. Abbeel,  
783 and J. Ba. Benchmarking model-based reinforcement learning. *CoRR*, abs/1907.02057,  
784 2019. URL <http://arxiv.org/abs/1907.02057>.
- 785 G. Wayne, C. Hung, D. Amos, M. Mirza, A. Ahuja, A. Grabska-Barwinska, J. Rae,  
786 P. Mirowski, J. Leibo, A. Santoro, M. Gemici, M. Reynolds, T. Harley, J. Abramson,  
787 S. Mohamed, D. Rezende, D. Saxton, A. Cain, C. Hillier, D. Silver, K. Kavukcuoglu,  
788 M. Botvinick, D. Hassabis, and T. Lillicrap. Unsupervised predictive memory in a  
789 goal-directed agent. *CoRR*, abs/1803.10760, 2018.
- 790 D. Yarats, A. Zhang, I. Kostrikov, B. Amos, J. Pineau, and R. Fergus. Improving sample  
791 efficiency in model-free reinforcement learning from images. In Q. Yang, K. Leyton-Brown,  
792 and Mausam (eds.), *Proceedings of the Thirty-Fifth Conference on Artificial Intelligence*  
793 *(AAAI’21)*. Association for the Advancement of Artificial Intelligence, AAAI Press, 2021.
- 794 A. Zhang, Z. Lipton, L. Pineda, K. Azizzadenesheli, A. Anandkumar, L. Itti, J. Pineau, and  
795 T. Furlanello. Learning causal state representations of partially observable environments.  
796 *CoRR*, 2019. URL <http://arxiv.org/abs/1906.10437>.
- 797 A. Zhang, S. Sodhani, K. Khetarpal, and J. Pineau. Learning robust state abstractions  
798 for hidden-parameter block MDPs. In *Proceedings of the International Conference on*  
799 *Learning Representations (ICLR’21)*, 2021. Published online: [iclr.cc](https://iclr.cc).
- 800 Y. Zhao, W. Zhao, R. Boney, J. Kannala, and J. Pajarinen. Simplified temporal consistency  
801 reinforcement learning. In *icml23*, 2023.
- 802 R. Zheng, X. Wang, Y. Sun, S. Ma, J. Zhao, H. Xu, H. Daumé, and F. Huang. Taco:  
803 Temporal latent action-driven contrastive loss for visual reinforcement learning. *Advances*  
804 *in Neural Information Processing Systems*, 36, 2024.
- 805 J. Zhu, Y. Xia, L. Wu, J. Deng, W. Zhou, T. Qin, T. Liu, and H. Li. Masked contrastive  
806 representation learning for reinforcement learning. *IEEE Transactions on Pattern Analysis*  
807 *and Machine Intelligence*, 45(3):3421–3433, 2022.



## 810 A APPENDIX

### 811 A.1 PROOF SKETCH OF PROPOSITION 3.1

812 In this section, we provide a theoretical foundation for the generalization capability of  
813 our proposed method. We formalize the relationship between subtask similarity and the  
814 embeddings learned by the GNN-based model. We restate the proposition in detail below:

815 **Proposition A.1.** *Let  $h_1, h_2, \dots, h_n \in \mathcal{H}$  be histories sampled from individual subtasks  
816 at different time steps in a POMDP, and let  $o'_1, o'_2, \dots, o'_n \in \mathcal{O}$  be the corresponding next  
817 observations. Let  $\phi: \mathcal{H} \rightarrow \mathcal{Z}$  be a belief function mapping histories to embeddings  $z_i = \phi(h_i)$ .  
818 Assume that  $\phi$  is Lipschitz continuous; that is, there exists a constant  $L_\phi > 0$  such that for  
819 all  $i, j$ :*

$$820 \|z_i - z_j\| \leq L_\phi \cdot d_{\mathcal{H}}(h_i, h_j),$$

821 where  $d_{\mathcal{H}}: \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}_{\geq 0}$  is a distance metric on  $\mathcal{H}$ . Let  $f: \mathcal{Z}^n \rightarrow \mathcal{O}$  be a model that  
822 predicts an observation  $o'_{\text{pred}} = f(z_1, \dots, z_n)$ . Assume that  $f$  is Lipschitz continuous with  
823 constant  $L_f > 0$ .

824 Then,

$$825 \left( \max_{i,j} d_{\mathcal{H}}(h_i, h_j) \leq \delta \right) \implies \mathcal{L}(o'_{\text{pred}}, o'_i) \leq (L_f L_\phi n \delta + \epsilon_i)^2,$$

826 where  $\epsilon_i$  represents the inherent error due to model approximation or noise.

#### 827 Proof Sketch.

828 **Step 1: Lipschitz Continuity of  $\phi$ .** Since  $\phi$  is Lipschitz continuous:

$$829 \|z_i - z_j\| \leq L_\phi \cdot d_{\mathcal{H}}(h_i, h_j) \leq L_\phi \delta \quad \text{for all } i, j.$$

830 **Step 2: Bounding Differences in Embeddings.** The maximum distance between any  
831 pair of embeddings  $z_i, z_j$  is bounded:

$$832 \|z_i - z_j\| \leq L_\phi \delta.$$

833 **Step 3: Lipschitz Continuity of  $f$ .** Applying  $f$  to embeddings  $z_1, \dots, z_n$  and another set  
834  $z'_1, \dots, z'_n$  (which in this case are  $z_j$ , since embeddings are close):

$$835 \|f(z_1, \dots, z_n) - f(z'_1, \dots, z'_n)\| \leq L_f \sum_{k=1}^n \|z_k - z'_k\|.$$

836 Since  $\|z_k - z'_k\| \leq L_\phi \delta$ :

$$837 \|f(z_1, \dots, z_n) - f(z'_1, \dots, z'_n)\| \leq L_f L_\phi n \delta.$$

838 **Step 4: Relating to the True Observation.** Assuming  $o'_j = f(z_j, \dots, z_j) + \epsilon_j$ , where  $\epsilon_j$   
839 accounts for model approximation error or noise. Then, for any  $i$ :

$$840 \|o'_{\text{pred}} - o'_i\| \leq \|o'_{\text{pred}} - o'_j\| + \|o'_j - o'_i\|.$$

841 Since  $o'_{\text{pred}}$  is close to  $o'_j$  due to the bound from Step 3, and  $o'_j$  is close to  $o'_i$  if  $o'_i \approx o'_j$ .

842 *Justification:* The model  $f$  processes a batch of embeddings  $z_1, \dots, z_n$  to predict the next  
843 observation  $o'_{\text{pred}}$ . When we input identical embeddings  $z_j$  into  $f$ , i.e.,  $f(z_j, \dots, z_j)$ , the  
844 model effectively focuses on the information contained in  $z_j$  without interference from  
845 variations in other embeddings. Given that  $z_j$  represents the embedding of history  $h_j$ , it  
846 is reasonable to expect that  $f(z_j, \dots, z_j)$  approximates the true next observation  $o'_j$ , up to  
847 some approximation error  $\epsilon_j$ .

848 **Step 5: Bounding the Prediction Error.** Combining the above:

$$849 \|o'_{\text{pred}} - o'_i\| \leq L_f L_\phi n \delta + \epsilon_i,$$

850 where  $\epsilon_i$  accounts for discrepancies between  $o'_i$  and  $o'_j$  and any inherent noise.



864 **Step 6: Squared Error Loss.** Therefore:

$$865 \mathcal{L}(o'_{\text{pred}}, o'_i) = \|o'_{\text{pred}} - o'_i\|^2 \leq (L_f L_\phi n \delta + \epsilon_i)^2.$$

866 Hence, minimizing the squared error loss under the Lipschitz continuity of  $\phi$  and  $f$  under the  
867 assumption of similar histories ensures that small differences in histories lead to proportionally  
868 small prediction errors. This confirms that our method effectively leverages relational  
869 structures among histories to generalize across subtasks, validating the proposition.  $\square$

870 While the proof establishes an upper bound on the prediction error based on the Lipschitz  
871 continuity of  $\phi$  and  $f$ , it's important to consider how minimizing the squared error loss

$$872 \mathcal{L}(o'_{\text{pred}}, o'_i) = \|o'_{\text{pred}} - o'_i\|^2$$

873 during training impacts the approximation errors  $\epsilon_i$  and the bound.

874 Minimizing  $\mathcal{L}$  reduces the approximation errors  $\epsilon_i$ , leading to a tighter bound on the prediction  
875 error:

$$876 \mathcal{L}(o'_{\text{pred}}, o'_i) \leq (L_f L_\phi n \delta + \epsilon_i)^2.$$

877 As  $\epsilon_i$  decreases, the bound becomes tighter, enhancing the model's predictive accuracy. This  
878 process improves the model's ability to generalize across similar histories and subtasks by  
879 effectively capturing relational structures in the data. Therefore, minimizing the loss during  
880 training is crucial for achieving the theoretical benefits outlined in the proof.

## 881 A.2 HYPERPARAMETERS AND EXPERIMENTAL DETAILS

882 Hyperparameter	883 Value
884 Discount factor ( $\gamma$ )	0.99
885 Number of environment steps	$3 \times 10^6$
886 Maximum number of distractors	4
887 Maximum size change	$12 \times 12$
888 Target network update rate ( $\tau$ )	0.005
889 Replay buffer size	400,000
890 Batch size	256
891 Learning rate	0.001
892 Latent state dimension	128
893 Epsilon greedy schedule	exponential(1.0, 0.05, 400,000)
894 R2D2 sequence length	10
895 R2D2 burn-in sequence length	5
896 $n$ -step TD	5
897 Training frequency	every 10 environment steps
898 Auxiliary loss coefficient ( $\lambda$ )	0.01
899 Latent state size	147
900 Num. neighbors in GNN ( $m$ )	4
901 Num. of message passing steps	2
902 Hidden state of Graph model	$147/2 = 73.5$

### 918 A.3 LATENT SPACE TRAJECTORIES

919  
920 This section outlines the methodology used to construct visual trajectories in the latent  
921 space of the encoder. These visualizations provide insights into how the latent spaces  
922 encode task-relevant information across different phases of the agent’s trajectory, such as  
923 key collection and goal navigation.

924 To generate these trajectories, we used the checkpoint of a trained encoder and simulated a  
925 path to the goal. We then divided this into two phases based on the subtask of key collection:  
926

- 927 1. **Phase 1:** trajectory until collection of the key.
- 928
- 929 2. **Phase 2:** trajectory after collecting the key until the goal.
- 930
- 931

932 For each phase, the hidden states produced by the encoder were collected during the execution  
933 of the corresponding actions. We then applied Principal Component Analysis (PCA) to  
934 reduce the dimensionality of these latent states to three components, enabling visualization  
935 in 3D space. The resulting points connect consecutive latent states, forming a trajectory  
936 in the latent space. Each connection and corresponding point is color-coded by phase to  
937 emphasize transitions between sub-tasks, with the goal state represented as a distinct point  
938 in the latent space. This visualization allows a qualitative comparison of how algorithms  
939 organize and structure their latent representations for task completion. We now summarize  
940 the general observations from these figures.

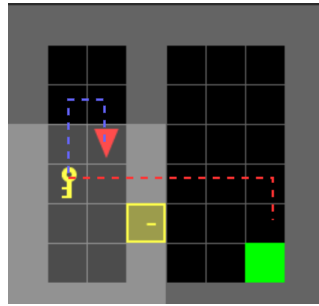
941  
942 **Clearer Trajectories in Graph\_OP.** The latent trajectories reveal notable differences  
943 in how various objectives shape the latent space representations. The **Graph\_OP** method  
944 consistently exhibits clearer and smoother trajectories between task phases, such as key  
945 collection and goal navigation. This clarity arises from the graph prediction objective, which  
946 helps the model learn a well-structured latent space. By focusing on observation prediction,  
947 **Graph\_OP** emphasizes encoding the environment’s dynamics and transitions between states,  
948 resulting in smoother and more structured latent.

949  
950 **Ruggedness in Graph\_AIS.** In contrast, incorporating the reward prediction objective, as  
951 seen in **Graph\_AIS**, introduces more ruggedness into the latent trajectories. This ruggedness  
952 reflects the aggressive influence of the reward prediction objective, which aligns the latent  
953 space with task rewards. While this alignment prioritizes encoding goal-directed information,  
954 it often disrupts the smooth structure typically learned by the graph prediction objective.  
955 Consequently, the latent trajectories for **Graph\_AIS** are less structured than that of **Graph\_OP**  
956 but better aligned with task-relevant rewards.  
957

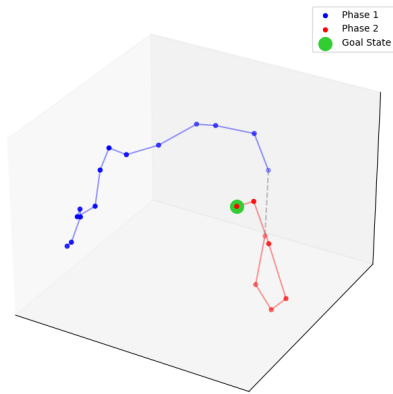
958  
959 **Goal State Placement.** Another key observation is the placement of the goal state in  
960 the latent space. In **Graph\_AIS**, the goal state appears further away from other latent states  
961 compared to **Graph\_OP**. This distinction highlights how the reward prediction objective drives  
962 the model to strongly differentiate goal states from other regions of the latent space. This  
963 explicit separation facilitates more effective credit assignment, enabling the agent to focus  
964 on actions that lead to the goal.  
965

966  
967 **Why Graph\_AIS Outperforms Graph\_OP.** Despite the less structured latent space,  
968 **Graph\_AIS** generally outperforms **Graph\_OP**. This is because reward alignment ensures that  
969 the latent space emphasizes task-relevant features, particularly those associated with long-  
970 term planning and goal achievement. Combining the graph and reward prediction objectives  
971 enables **Graph\_AIS** to balance relational modeling and goal-directed alignment, improving  
task performance.

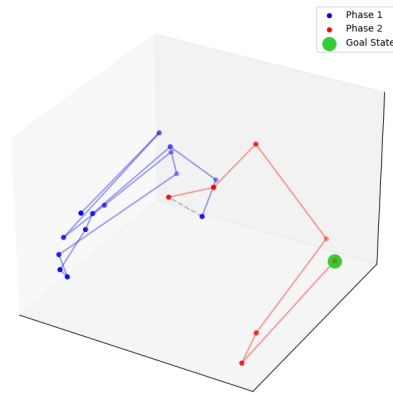
A.3.1 MINIGRID-DOORKEY-8x8-v0



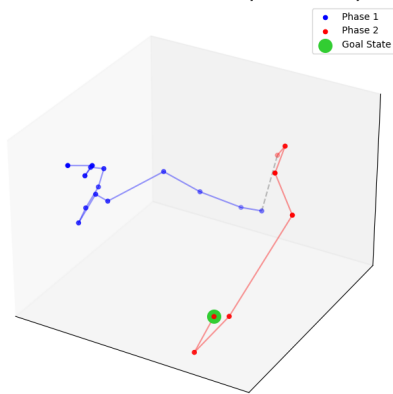
(a) Simulated Trajectory



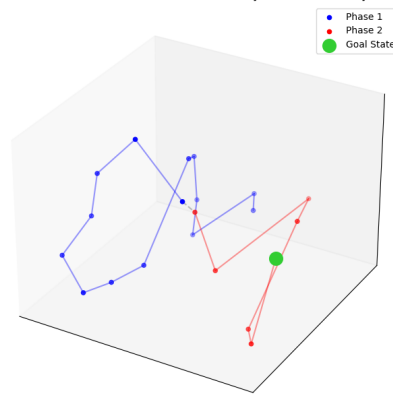
(b) Graph\_OP



(c) Graph\_AIS

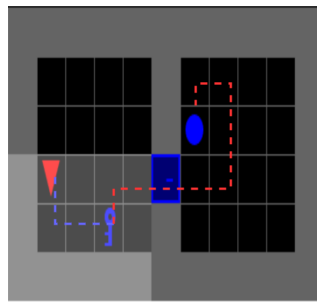


(d) min\_OP

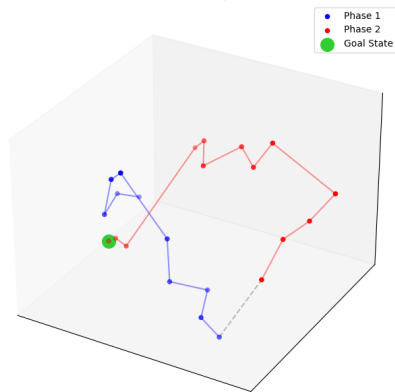


(e) min\_AIS

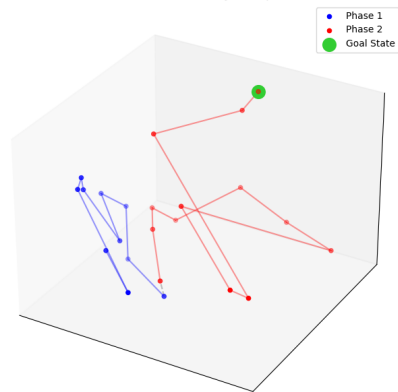
1026 A.3.2 MINIGRID-OBSTRUCTEDMAZE-1DL-V0  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048



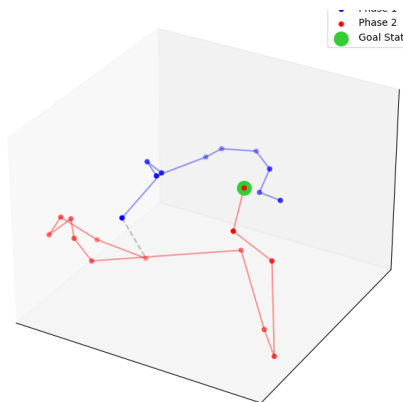
1049 (a) Simulated Trajectory



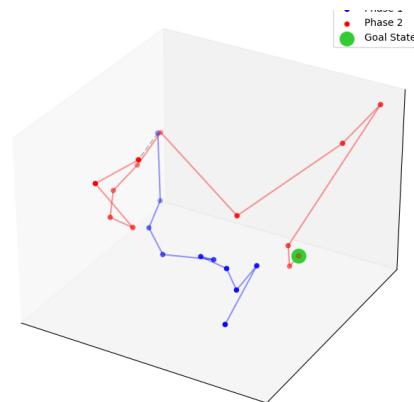
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063 (b) Graph\_OP



1064 (c) Graph\_AIS

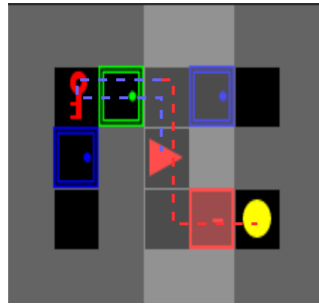


1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079 (d) min\_OP

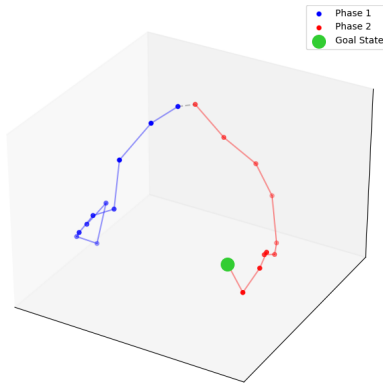


(e) min\_AIS

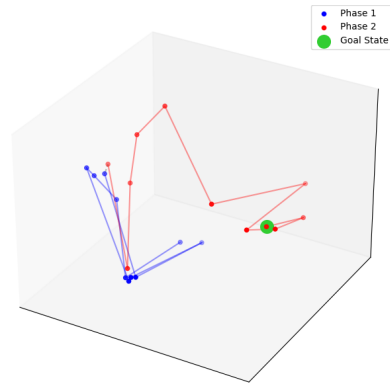
1080 A.3.3 `MINIGRID-KEYCORRIDORS3R2-V0`  
 1081  
 1082  
 1083  
 1084  
 1085  
 1086  
 1087  
 1088  
 1089  
 1090  
 1091  
 1092  
 1093



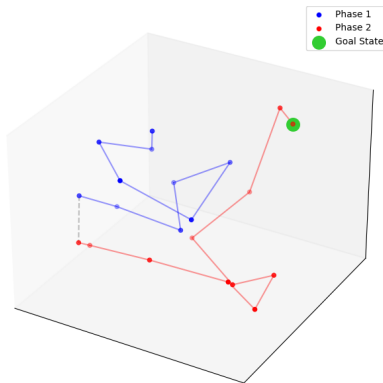
(a) Simulated Trajectory



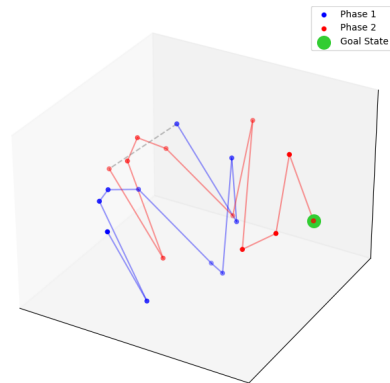
(b) Graph\_OP



(c) Graph\_AIS



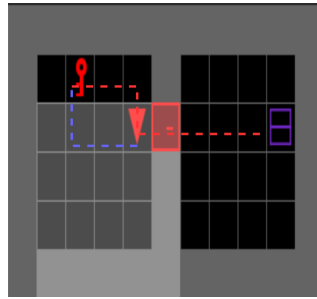
(d) min\_OP



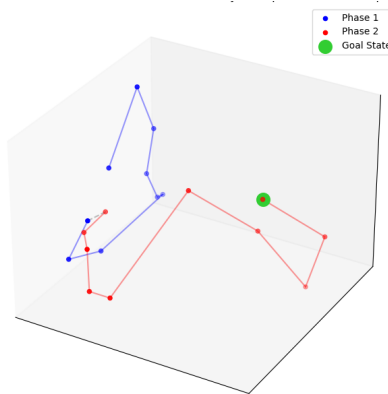
(e) min\_AIS

1130  
 1131  
 1132  
 1133

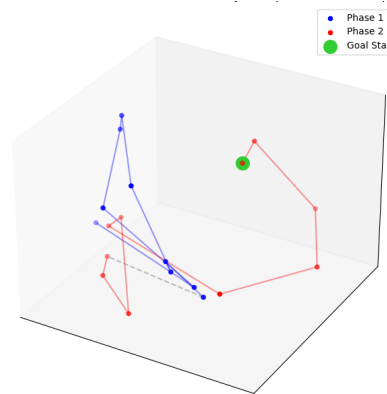
1134 A.3.4 MINIGRID-UNLOCKPICUP-V0  
1135



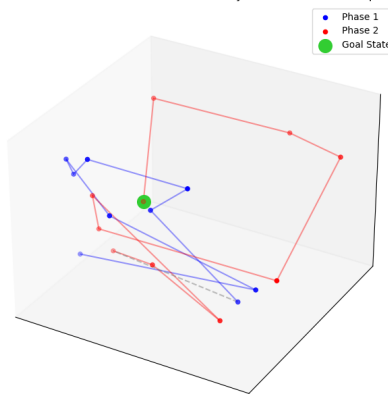
(a) Simulated Trajectory



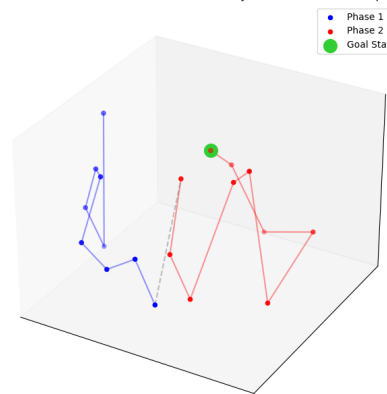
(b) Graph\_OP



(c) Graph\_AIS



(d) min\_OP



(e) min\_AIS

1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187



#### 1188 A.4 PREDICTION OF THE MODELS

1189

1190

1191 In this section, we compare the predictions generated by the MLP-based model and the  
 1192 Graph-based model. To produce the predictions in the figures below, we initialized an agent,  
 1193 loaded the model, critic, and encoder checkpoints, and populated the buffer by interacting  
 1194 with the environment. A minibatch of observations was then sampled from this buffer, and  
 1195 the model was queried to predict the corresponding subsequent observations. The figure  
 1196 compares an observation image from the batch with the predictions from the MLP-based  
 1197 and Graph-based models.

1198 The Graph-based model consistently generates predictions with higher fidelity than the  
 1199 MLP-based model, highlighting the advantages of the GNN’s temporal reasoning capabilities.  
 1200 While the MLP model struggles to produce visually accurate reconstructions, it retains  
 1201 vital features such as approximate spatial contrasts and object colors. These features may  
 1202 explain its ability to perform reasonably despite poor visual quality. In contrast, the Graph  
 1203 model produces predictions that closely resemble the original observations, demonstrating  
 1204 its superior ability to leverage temporal relationships across trajectories.

1205

1206

1207

##### 1208 A.4.1 MINIGRID-DOORKEY-8x8-v0

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

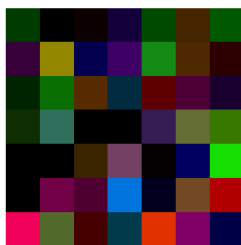
1219

1220

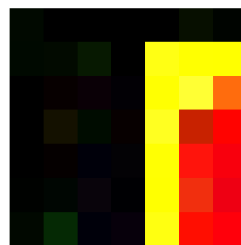
1221



1222 (a) Observation



1223 (b) MLP Model



1224 (c) Graph Model

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

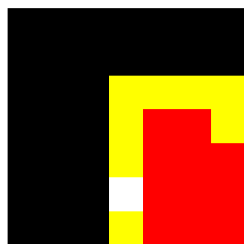
1237

1238

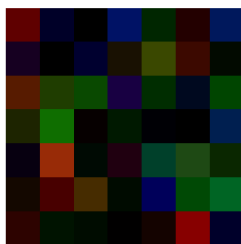
1239

1240

1241



1242 (a) Observation



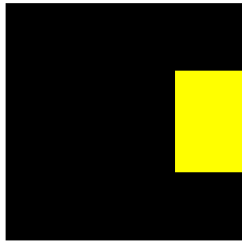
1243 (b) MLP Model



1244 (c) Graph Model

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

A.4.3 `MINIGRID-KEYCORRIDORS3R2-V0`



(a) Observation



(b) MLP Model

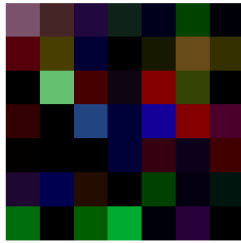


(c) Graph Model

A.4.4 `MINIGRID-UNLOCKPICKUP-V0`



(a) Observation



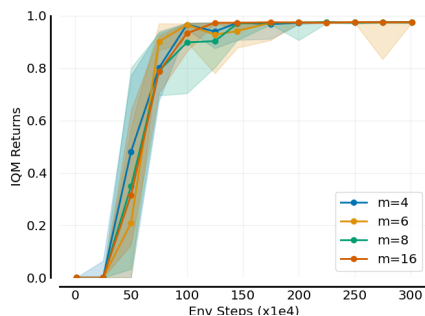
(b) MLP Model



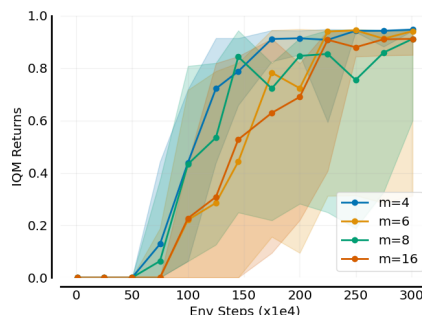
(c) Graph Model

### A.5 DIFFERENT VALUES FOR NEIGHBORS

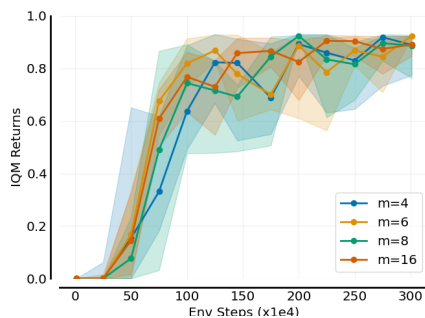
We ablated the number of neighbors ( $m$ ) used in the graph construction to evaluate its effect on task performance. The results, presented in Appendix A.5, demonstrate that the model is robust to changes in  $m$ , with similar final returns across  $m = 4$ ,  $m = 6$ ,  $m = 8$ , and  $m = 16$  in most tasks. In the early stages of training,  $m = 4$  tends to achieve faster returns, suggesting that smaller graphs may provide more efficient learning initially. However, tasks with more complex relational dependencies, such as `UnlockPickup-v0`, benefit slightly from  $m = 6$ , indicating that the optimal number of neighbors may be task-specific. Larger values of  $m$  introduce more variability in performance for some environments, as evidenced by broader confidence intervals, potentially due to increased noise in the graph. Overall, these results highlight the robustness of the proposed method across different graph configurations, with  $m = 4$  serving as a reasonable default choice for most tasks.



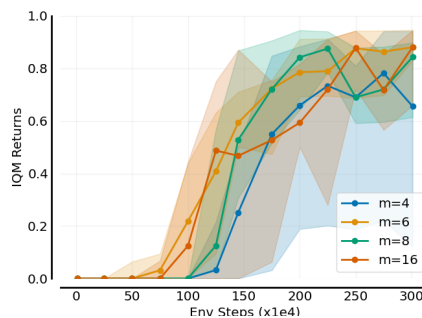
(a) MiniGrid-DoorKey-8x8-v0



(b) MiniGrid-ObstructedMaze-1D1-v0



(c) MiniGrid-KeyCorridorS3R2-v0



(d) MiniGrid-UnlockPickup-v0

### A.6 ISOLATING THE EFFECT OF RELATIONAL REASONING IN THE GNN

We designed an experiment to isolate this effect and understand whether the GNN’s observed benefits arise from its relational reasoning or simply from operating on the entire batch of observations. In our standard setup, the GNN processes a batch of observations by constructing a graph over the entire batch and performing relational reasoning through message passing. By contrast, the baseline MLP independently predicts the next observation for each element in the batch without leveraging relationships across the batch.

We modified the GNN and MLP architectures for this experiment to process mini-batches of 50 observations each sequentially. Specifically, we divided the original batch into 50-unit mini-batches and processed them sequentially. The GNN constructed a graph over each mini-batch and performed relational reasoning with a sparse connection via message passing, while the MLP processed the mini-batches without relational reasoning. After processing each mini-batch, the outputs were concatenated into a new batch with the same dimensionality as the original input, and a final linear transformation was applied to produce the output.

This setup ensures that both architectures operate sequentially on mini-batches, making the primary difference between them using relational reasoning in the GNN. The results, shown in Figure 15, demonstrate that the GNN-based model outperforms the MLP-based model in this scenario, indicating that the benefits of the GNN arise from its ability to reason over observations within each mini-batch relationally. This experiment highlights the critical role of relational reasoning in achieving better performance.

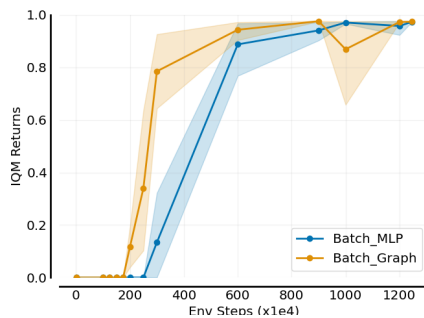


Figure 15: Difference between batch of 50 observations for the Graph and MLP models MiniGrid-UnlockPickup-v0