

# Towards better dense rewards in Reinforcement Learning Applications

Shuyuan Zhang<sup>1,2</sup>

<sup>1</sup>School of Computer Science, McGill University

<sup>2</sup>Mila

## 1 Introduction

Finding meaningful and accurate dense rewards is a fundamental task in the field of reinforcement learning (RL) that enables agents to explore environments more efficiently. In traditional RL settings, agents learn optimal policies through interactions with an environment guided by reward signals. However, when these signals are sparse, delayed, or poorly aligned with the intended task objectives, agents often struggle to learn effectively. Dense reward functions, which provide informative feedback at every step or state transition, offer a potential solution by shaping agent behavior and accelerating learning.

Despite their benefits, poorly crafted reward functions can lead to unintended behaviors, reward hacking, or inefficient exploration. This problem is particularly acute in complex or high-dimensional environments where handcrafted rewards are difficult to specify and validate.

To address this, recent research has explored a variety of approaches, including inverse reinforcement learning, reward modeling from human preferences, and self-supervised learning of intrinsic rewards. While these methods offer promising directions, they often involve trade-offs between generality, scalability, and alignment with human intent. This proposal explores several approaches to dealing with these unsolved problems and enhancing the effectiveness and reliability of dense reward construction in different RL applications.

### Proposal Structure

- Section 2 introduces the relevant notation and background knowledge.
- In Section 3, we present Graph-Guided subGoal representation Generation (G4RL), a novel graph encoder-decoder framework for goal-conditioned scenarios. This work’s primary objective is to boost the performance of existing Goal-Conditioned Hierarchical Reinforcement Learning (GCHRL) methods. Traditional GCHRL often relies on the Euclidean distance between the current state and the assigned subgoal to generate an intermediate reward for the low-level agent. This reliance can be problematic, as the geometric distance in the state space often fails to align with the true transition distance or accessibility between states. Consequently, this leads to inaccurate and misleading intermediate reward signals. Our G4RL framework addresses this by introducing a method to encode subgoal representations using the underlying transition graph of the environment. By leveraging this graph, our encoder-decoder architecture generates subgoal embeddings that inherently reflect the true transition dynamics. This approach produces more accurate and informative dense rewards, significantly enhancing the effectiveness of typical GCHRL methods.
- Section 4 addresses the challenge of sparse rewards in traditional Reinforcement Learning from Human Feedback (RLHF) and introduces our proposed solution: Shapley Credit Assignment Rewards (SCAR). Vanilla RLHF suffers from reward sparsity because the reward model typically provides a usable signal only upon the completion of a full response. This lack of intermediate feedback hinders efficient learning. While previous methods have attempted to use dense rewards by evaluating partial sequences, they often introduce inaccurate rewards because the partial sequences lie outside the reward model’s training distribution (out-of-distribution). Our methodology resolves both issues simultaneously. It leverages the Shapley value from cooperative game theory to assign a statistically sound credit score to each

component (e.g., token or phrase) of a full generated sentence. This approach effectively converts the sparse, sentence-level reward into an accurate, dense signal for every part of the sequence, thus promoting more robust and efficient learning.

- In Section 5, we outline an ongoing project focused on enhancing agents’ continual learning capabilities. This is achieved by introducing an auxiliary reward signal derived from structural similarities among tasks. The core idea is to guide the agent’s learning process by having it identify structural commonalities across different tasks. When the agent encounters a new task whose underlying structure is similar to previous experiences, it leverages this recognized similarity as a crucial learning heuristic and guidance signal, thereby accelerating adaptation and improving performance in the continual learning setting.

## 2 Background

### 2.1 Markov Decision Processes

As the most common framework for modeling reinforcement learning scenarios, the Markov Decision Process (MDP) [39] is introduced as a tuple  $\langle \mathcal{S}, \mathcal{A}, P, R, \gamma \rangle$ , defined as follows: At each time step  $t$ , the agent observes the current state  $s_t \in \mathcal{S}$  provided by the environment and chooses an action  $a_t \in \mathcal{A}$  according to its internal policy  $\pi(a_t|s_t)$ , which specifies the probability of choosing action  $a_t$  given state  $s_t$ . The action is then executed, and the interaction with the environment leads the agent to a new state  $s_{t+1}$  according to a transition probability function  $P(s_{t+1}|s_t, a_t)$  which is known only to the environment. Subsequently, the agent receives a reward  $r_t$ , determined by the reward function  $R(s_t, a_t)$  that evaluates the action taken in the current state and is also only visible to the environment. The agent aims to learn an optimal policy  $\pi$  to maximize the expected discounted cumulative reward  $\mathbb{E}_\pi[\sum_{t=0}^T \gamma^t r_t]$ , where  $\gamma$  (with  $0 \leq \gamma < 1$ ) is a pre-defined discount factor used to prioritize immediate rewards over distant future rewards, thereby ensuring that the total reward remains finite.

### 2.2 Goal-conditioned Hierarchical RL (GCHRL)

Goal-Conditioned Reinforcement Learning (GCRL) trains agents to achieve specific goals, which are the target states. The agent receives an additional goal input  $g_t$  along with the state input  $s_t$  and learns a policy  $\pi(a_t|g_t, s_t)$  that aims to achieve this goal. Goals are represented explicitly in the input to the policy, guiding the agent’s actions towards desired outcomes. The reward function is goal-dependent, providing positive feedback when the agent successfully reaches the desired goal state.

To deal with large and complex environments, Goal-Conditioned Hierarchical Reinforcement Learning (GCHRL) [33, 52] decomposes the learning task into a hierarchy of smaller, more manageable sub-tasks. Typically, there are two levels of agents. At time step  $t$ , the high-level agent chooses a subgoal  $g_t$ , a representation of a target state, and assigns it to the low-level agent to achieve as part of the overall task. This choice is made by sampling  $g_t$  from the high-level policy  $\pi_h(g_t|\phi(s_t))$ , where  $\phi : s \mapsto \mathbb{R}^d$  is the state representation function which gives a condensed representation of the state.

Each state  $s_t$  can be mapped to its subgoal feature  $g(s_t)$  by a subgoal feature extractor. Note that  $g(s_t)$  is not the same as  $g_t$ . The former,  $g(s_t)$ , is the learned subgoal feature of the current state  $s_t$ , while the latter,  $g_t$ , is the target state we aim to reach from the state  $s_t$  in one step or multiple steps.

Given the subgoal  $g_t$  sampled from the high-level policy  $\pi_h(g_t|\phi(s_t))$  for the current time step  $t$  and the state representation vector  $\phi(s_t)$ , a low-level agent executes action  $a_t$  based on the low-level policy  $\pi_l(a_t|\phi(s_t), g_t)$ . The low-level agent is trained using the intrinsic (dense) reward signal  $r_{\text{int}}(s_t, g_t, a_t, s_{t+1}) = -\|\phi(s_{t+1}) - g_t\|_2^2$  to encourage it to achieve the subgoal.

Both agents can be implemented by any policy-based methods, including those introduced in previous works on policy gradients such as [19, 21] and [42].

### 2.3 RL for Text Generation

Reinforcement Learning (RL) has been increasingly leveraged for fine-tuning LLMs in text generation tasks [41, 35, 27, 10, 6]. Unlike standard supervised fine-tuning which relies on maximizing the likelihood of ground-truth sequences, RL enables optimization directly towards sequence-level objectives or metrics that are not differentiable, such as ROUGE scores in summarization or human preferences for dialogue quality [44, 37]. A particularly successful application of this is Reinforcement Learning from Human Feedback (RLHF)

[37, 14, 44, 7], which has become a widely adopted technique for aligning LLMs with complex human values and instructions. The typical RLHF process involves learning a reward model (RM) from a dataset of human comparisons between different model outputs, followed by fine-tuning the LLM policy to maximize the scalar reward assigned by the RM to the generated text sequences. However, a widely acknowledged challenge in this standard RLHF framework is the inherent sparsity of the reward signal: the RM typically provides feedback only after the entire sequence has been generated [44, 7]. This terminal reward makes the temporal credit assignment problem—identifying which specific token choices (actions) contributed positively or negatively to the final outcome—particularly difficult [45]. This difficulty has spurred research into methods for creating denser and more informative reward signals.

### 3 Using Spatial Information as Dense Reward in Goal-Conditioned Hierarchical Reinforcement Learning

This section presents our first published work, Graph-Guided subGoal representation Generation (G4RL) [51]. This method reshapes the subgoal space by utilizing a state graph to incorporate the relative spatial information of visited states, thereby generating better dense rewards for different levels of agents in goal-conditioned RL.

One drawback of previous hierarchical reinforcement learning algorithms [33, 24, 52, 30] is that the Euclidean distance calculated in the original state representation space between the current state and the intended goal does not accurately reflect the true progress of the low-level agent, as there is rarely a straight path between the current state and the subgoal in the space.

As a result, the low-level agent trained with such information may receive an inaccurate reward signal, thus impairing its performance. Another issue is that, without appropriate constraints, the high-level agent may propose subgoals that are too difficult to reach, wasting exploration steps on pursuing infeasible targets [52]. The proposed method aims to mitigate both problems by calculating the distance in a subgoal representation space between subgoal representations given by a graph encoder-decoder. This graph encoder-decoder captures the actual connectivity between states, ensuring that the generated subgoal representations respect adjacency information.

#### 3.1 State graph

To record the visited states and their connections, We maintain a state graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  with a fixed number  $N$  of nodes<sup>1</sup>. This graph is built and updated during training, with no pre-training using expert data or handcrafted process involved in its construction.

Each node is labelled by the corresponding state and for each node  $s_t$ , the corresponding state representation vector  $\phi(s_t)$ , which is also referred to as the node feature, is stored. Edges in the graph represent connectivity between states. The graph is constantly updated during exploration.

We choose the state graph to be undirected for mainly two reasons: **(1) Efficiency:** An undirected graph requires fewer resources for storage and computation compared to a directed graph, and **(2) Method compatibility:** The method relies on defining a deterministic distance between state/subgoal representations for each node pair. This is straightforward in an undirected graph but becomes problematic in a directed setting, where distances can be asymmetric or undefined.

This choice is based on the assumption that G4RL is designed for environments with symmetrical and reversible dynamics. However, some of the experiments demonstrate that G4RL can still enhance performance in partially asymmetric environments; please refer to the experiments section for further details.

##### 3.1.1 Graph construction

The graph is initialized with  $N$  empty nodes and no edges. The corresponding weighted adjacency matrix  $\mathbf{A}$  is set to an  $N \times N$  zero matrix. We perform the GCHRL exploration process using randomly initialized policy  $\pi_h$  and  $\pi_l$ . Once the agent encounters a state representation never seen before, that is, the representation is different from any state representations stored in the graph, as described in equation (1), it stores the state

<sup>1</sup>The number of training states for the graph encoder-decoder grows quadratically with  $N$  because the adjacency weight matrix has  $N^2$  elements. The choice of  $N$  depends on the machine’s capabilities.

representation  $\phi(s_t)$  as the node feature of an empty node in the graph and build an edge between this node and the node corresponds to the previous state:

$$\forall_{s_v \in \mathcal{V}}, \|\phi(s_t) - \phi(s_v)\|_2 > \epsilon_d, \quad (1)$$

$$\mathbf{A}_{\phi(s_t), \phi(s_{t-1})} = \mathbf{A}_{\phi(s_{t-1}), \phi(s_t)} = 1, \quad (2)$$

where  $\epsilon_d$  is a hyper-parameter controlling the distance threshold between state representations.

When the agent encounters a state  $s_t$  with feature  $\phi(s_t)$  that is similar to several node representations already stored in the graph, it finds the state whose representation is the closest to the current state feature:

$$s_v = \arg \min_{s_u: \|\phi(s_t) - \phi(s_u)\|_2 \leq \epsilon_d} \|\phi(s_t) - \phi(s_u)\|_2. \quad (3)$$

Then the node  $s_v$  is relabeled as  $s_t$  and the weight for the edge  $(s_{t-1}, s_t)$  is updated as follows:

$$\mathbf{A}_{\phi(s_{t-1}), \phi(s_t)} = \mathbf{A}_{\phi(s_t), \phi(s_{t-1})} := \mathbf{A}_{\phi(s_{t-1}), \phi(s_t)} + 1. \quad (4)$$

Note that a large weight indicates more frequent transitions between the underlying states.

We have used the Euclidean norm to define the distance between feature vectors. Since some elements may contain more spatial information than others, one can use a weighted Euclidean norm to define the distance between state representations instead.

### 3.1.2 Graph updating

The graph has a fixed number of nodes. Suppose the graph is now full. When a new state  $s_t$  is encountered, if  $s_v$  from equation (3) exists, as before we relabel the node as  $s_t$  and perform edge update as shown in equation (4); Otherwise, we replace the oldest state node in the graph with the current state node, delete all edges previously linked to that node, and create an edge  $(\phi(s_{t-1}), \phi(s_t))$  with weight  $\mathbf{A}_{\phi(s_{t-1}), \phi(s_t)} = \mathbf{A}_{\phi(s_t), \phi(s_{t-1})} = 1$ . Alternatively, we could replace the state node that is most weakly connected to the other nodes—that is, the node with the lowest sum of edge weights.

## 3.2 Graph encoder-decoder

To enable the assignment of suitable subgoal representations to every possible state, including unseen ones, we use node representations and edges to train a graph encoder-decoder. The parameter updates of the graph encoder-decoder and the policies during policy training are performed alternately in each episode.

The encoder-decoder starts training after the graph is full and continues periodically after processing a few trajectories. Section 3.3 will show the details of the training schedule.

The encoder  $\mathbf{E}$  maps every state representation  $\phi(s)$  to a subgoal representation  $g(s)$ . We use a feed-forward network (FFN) with several layers as the encoder  $\mathbf{E}$ :

$$g(s) = \mathbf{E}(\phi(s)) = \text{FFN}(\phi(s)). \quad (5)$$

The weight parameters of the feed-forward network will be learned through training. The decoder  $\mathbf{D}$  takes two subgoal representations as input and outputs the inner product of these two representations:

$$\mathbf{D}(g(s_u), g(s_v)) = g(s_u)^T g(s_v). \quad (6)$$

We choose dot-product similarity based on the assumption that the similarity between two nodes, such as the overlap in their local neighbourhoods, is well captured by the dot product of their embeddings. This assumption is supported by prior work in the graph embedding literature [2, 12, 36].

The aim is to use the encoder-decoder structure to predict node relations. Naturally we can use  $\mathbf{A}_{\phi(s_u), \phi(s_v)}$  as a measure of the relation between nodes  $\phi(s_u)$  and  $\phi(s_v)$ . But for the sake of numerical stability in the training process, we use  $\mathbf{A}_{\phi(s_u), \phi(s_v)} / \max_{\phi(s_i), \phi(s_j)} \mathbf{A}_{\phi(s_i), \phi(s_j)}$  as a measure.

Thus the loss function is defined as:

$$\mathcal{L} = \sum_{\phi(s_u), \phi(s_v) \in \mathcal{V}} [\mathbf{D}(g(s_u), g(s_v)) - \mathbf{A}_{\phi(s_u), \phi(s_v)} / \max_{\phi(s_i), \phi(s_j)} \mathbf{A}_{\phi(s_i), \phi(s_j)}]^2. \quad (7)$$

This loss function can enforce the subgoal representation provided by the encoder to respect neighbouring features in the graph.

Note that in each training phase of the graph encoder-decoder (except the first one), we use the values of the parameters obtained from the last training phase as the initial point, which helps save computation cost.

### 3.3 Adaptive training schedule of the graph encoder-decoder

The graph stores evolving data, including state representations as node features and connection information in the weighted adjacency matrix  $\mathbf{A}$ , which are continuously updated during online training. Since the graph structure and content change at varying rates across episodes, training the graph encoder-decoder at fixed intervals can cause several issues: (1) high variance in earlier episodes, where sparse or unstable graph data may lead to unreliable model updates; (2) data underutilization, where intermediate graph states are overwritten before being used for training; and (3) overfitting in later episodes, as the model repeatedly trains on increasingly redundant data. To address these issues, we introduce an adaptive training schedule for the graph encoder-decoder, described in the following paragraph.

There are two types of data changing in the graph: node replacement and edge update. We introduce a variable  $c$  to track the weighted number of data changes. Since the replacement of nodes has a much higher impact on the data than the edge update, we add  $N - 1$  to  $c$  if a node replacement occurs, and add 1 to  $c$  if an edge update happens:

$$c = \begin{cases} c + N - 1, & \text{if a node replacement happens,} \\ c + 1, & \text{if an edge update happens.} \end{cases} \quad (8)$$

When this variable exceeds a certain value, specifically a tolerance  $\beta$  multiplied by the total number of non-diagonal elements  $N^2 - N$  in the matrix  $\mathbf{A}$ , we perform one gradient update for the graph encoder-decoder, and then we reset  $c$  to 0.

### 3.4 Hierarchical agent with graph encoder-decoder

The proposed method involves traditional goal-conditioned settings and a subgoal representation extractor implemented by a graph encoder-decoder.

The high-level policy  $\pi_h(g_t|\phi(s_t))$  nominates a subgoal every  $K$  steps and is trained using the external environmental reward  $r_{\text{ext}}$ .

The policy can be implemented by any policy-based RL algorithm that takes transition tuples  $(s_t, g_t, a_t, r_t, s_{t+1}, g_{t+1})$  as input. To encourage it to propose a subgoal that is not too difficult to reach from the current state  $s_t$  for more efficient exploration, we add an intrinsic term to the high-level reward, considering the distance between the subgoal features of  $s_t$  and  $g_t$  in the subgoal space:

$$r_h(s_t, g_t, s_{t+1}) = r_{\text{ext}} + r_{\text{int}} = r_{\text{ext}} + \alpha_h \cdot \mathbf{D}(\mathbf{E}(\phi(s_t)), \mathbf{E}(g_t)), \quad (9)$$

where  $\alpha_h$  is a hyperparameter that controls the significance of the intrinsic term in the high-level reward.

The low-level policy  $\pi_l(a_t|\phi(s_t), g_t)$ ,

however, operates in the subgoal space. While it still takes  $\phi(s_t)$  and  $g_t$  as input and outputs an atomic action  $a_t$ , we compute the reward based on distances in the subgoal space:

$$r_l(s_t, g_t, a_t, s_{t+1}) = -\|\phi(s_{t+1}) - g_t\|^2 + \alpha_l \cdot \mathbf{D}(\mathbf{E}(\phi(s_{t+1})), \mathbf{E}(g_t)), \quad (10)$$

where  $\alpha_l$  is a hyperparameter controlling the significance of the reward term in the low-level reward. By computing the intrinsic reward in the subgoal space rather than in the state space, the function provides high values when proposed subgoals are easy to reach from the current location and low values when subgoals are close in the original state space but difficult to reach from the current location. The low-level agent can also be any policy-based algorithm.

### 3.5 Balancing between speed and performance

Due to the excessive comparisons between the current state representation and the node features during graph updates,

as well as the training cost of the graph encoder-decoder on a large graph, the experiments show that the GCHRL method, after incorporating the method, takes approximately twice as long as the original GCHRL method.

To reduce the additional cost, we can either decrease the frequency of sampling candidates for node features, train the graph encoder-decoder with a subset of all available training data, or do both.

For the sampling frequency, instead of comparing the state representation with node features in each time step, we do it in every  $t_c$  time steps. This may significantly speed up the method while maintaining satisfactory performance.

To reduce the training data in each graph encoder-decoder training cycle, we randomly sample node pairs in the graph instead of using every node pair for training.

In the experiments section, we will present the training time and performance of G4RL with the above two techniques applied.

### 3.6 Results

**Environmental settings** We used AntMaze, AntGather, AntPush, AntFall and AntMaze-Sparse environments from the GYM MuJoCo library [46]. The first four involve complex navigation and manipulation tasks performed by a simulated multi-armed robot, while AntMaze-Sparse presents a particularly challenging scenario due to sparse reward signals, providing feedback only upon reaching the goal. Note that AntPush and AntFall contain asymmetric transitions which lead to irreversible state changes that can be difficult to learn with an undirected graph. We added these environments to show G4RL’s robustness in inherently asymmetric environments. For the state representation, we selected a subset of raw state dimensions that contains spatial information (e.g. coordinates and arm angles) to serve as the node representation for the graph encoder-decoder. We deliberately aligned the choice of environments with prior work to ensure a fair and consistent comparison between the backbone algorithms and their G4RL-augmented versions. This allowed us to use the same environments, hyperparameters, and codebases provided by the original studies. Our goal was to demonstrate that G4RL can consistently enhance the performance of these backbone algorithms under comparable conditions.

**Baselines** For baseline methods, we select typical GCHRL methods that use a Euclidean distance as the intermediate reward for the low-level agent to assess G4RL’s capability to correct the bias of this reward and improve the performance of these baselines.

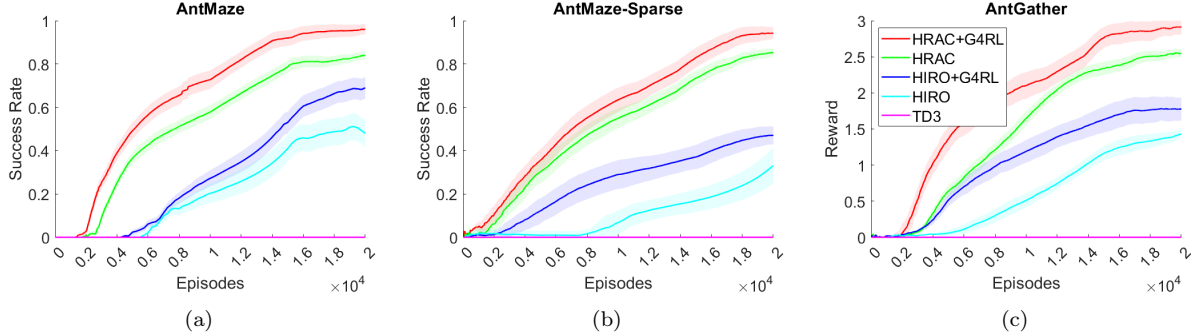


Figure 1: Success Rate on (a) AntMaze (b) AntMaze-Sparse and Reward on (c) AntGather, using HIRO, HIRO-G4RL, HRAC, HRAC-G4RL, and TD3. Incorporating G4RL in HIRO and HRAC significantly enhances their performance.

The learning curves of baseline methods and G4RL-applied versions are plotted in Figure 1 and 2. Note that all the curves reported are averages from 20 independent runs and they have been equally smoothed for better visualization.

From Figures 1 and 2, we observe that, in all environments, incorporating G4RL in the base GCHRL methods significantly enhances their performance, further improving the already strong results of these hierarchical methods compared to the non-hierarchical method. Notably, G4RL-augmented methods not only achieve higher final success rates but also converge substantially faster, with the most significant improvements observed during the early stages of training.

To demonstrate the proposed method’s effectiveness in environments with image-based state representations, we conducted experiments on AntMaze, AntPush, and AntFall, utilizing images as states, and compared the results with HESS and HLPS, along with their G4RL variations. We use Mean Squared Error (MSE) to measure the pixel-wise differences between image states to decide whether a new node should be added to the

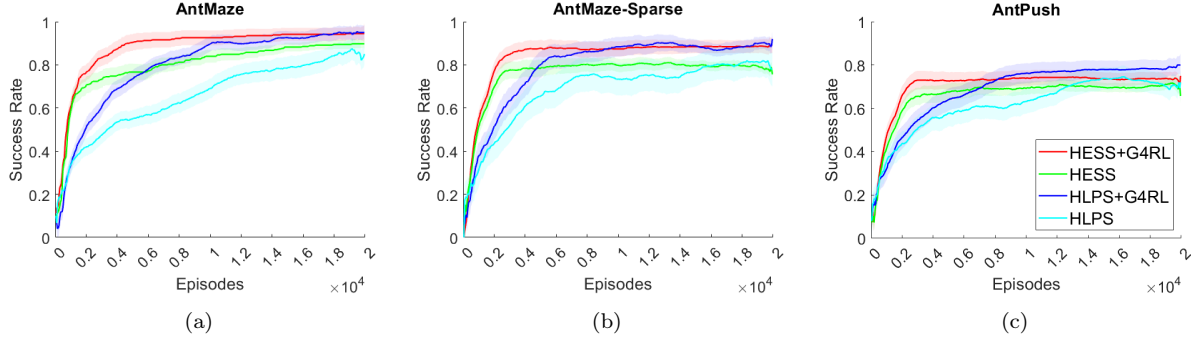


Figure 2: Success Rate on (a) AntMaze (b) AntMaze-Sparse and (c) AntPush, using HESS, HESS-G4RL, HLPS, HLPS-G4RL. Incorporating G4RL in HESS and HLPS significantly enhances their performance.

graph. The test results, given in Figure 3, show that methods incorporating G4RL exhibit faster convergence and achieve higher performance across all tested image-based environments.

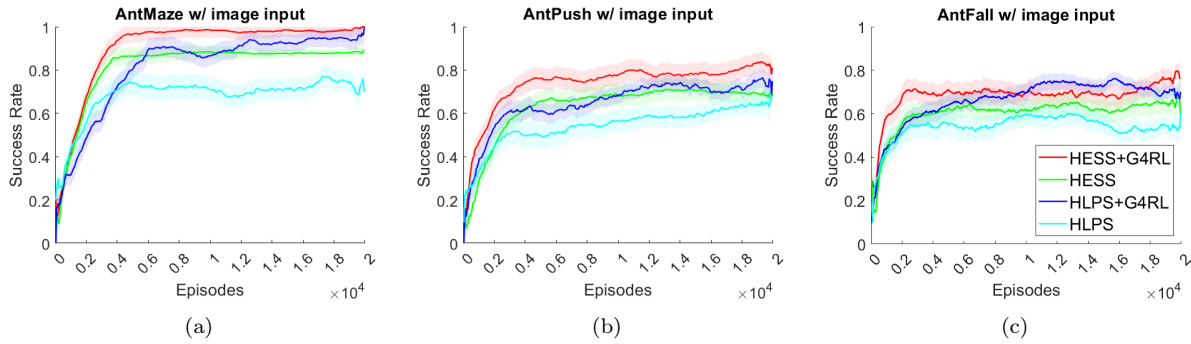


Figure 3: Success Rate on (a) AntMaze (b) AntPush and (c) AntFall with image state features, using HESS, HESS-G4RL, HLPS, HLPS-G4RL. Incorporating G4RL helps convergence and achieves higher performance across all tested image-based environments.

### 3.6.1 The effect of high/low-level intrinsic reward

We consider the following variants of G4RL to show the effectiveness of adding high-level and low-level intrinsic rewards:

- **High+Low-level intrinsics:** Apply both equation (9) and equation (10) to the high-level and low-level rewards respectively.
- **High-level intrinsic only:** Apply equation (9) to the high-level rewards and set  $\alpha_l = 0$  in equation (10) when it is applied to the low-level rewards.
- **Low-level intrinsic only:** Apply equation (10) to the low-level rewards and set  $\alpha_h = 0$  in equation (9) when it is applied to the high-level rewards.
- **HIRO/HRAC/HESS/HLPS:** Vanilla baseline methods.

Same as before, all the curves reported in this section are drawn from results averaged across 20 independent runs. All curves have been equally smoothed for better visualization.

Figures 4 to 5 show that, across all tested environments and algorithms, the combination of high-level and low-level intrinsic rewards results in the highest success rates and fastest convergence. The high-level intrinsic-only variant outperforms the low-level intrinsic-only variant, especially in sparse reward tasks, indicating that high-level intrinsic rewards play a crucial role in facilitating efficient exploration by encouraging the agent to select reachable and meaningful subgoals. In contrast, low-level intrinsic rewards have limited effect

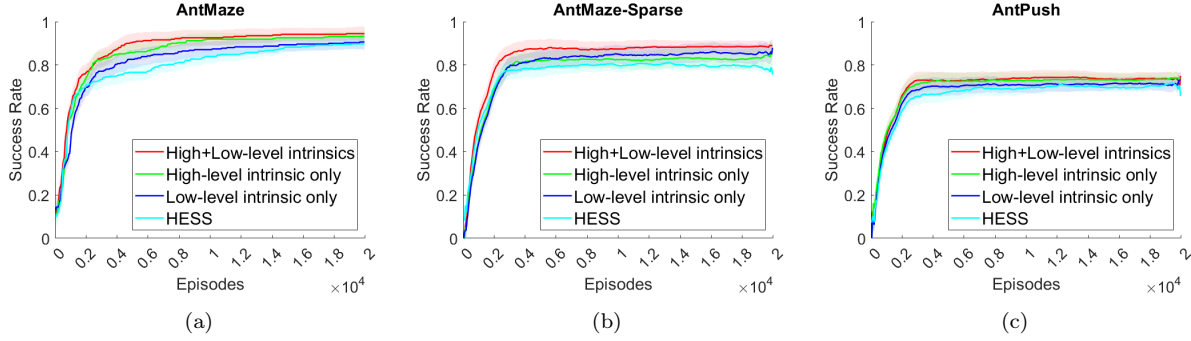


Figure 4: Success Rate on (a) AntMaze (b) AntMaze-Sparse and (c) AntPush using HESS-G4RL, HESS + High-level intrinsic, HESS + Low-level intrinsic and HESS. All curves have been equally smoothed for better visualization. The combination of high-level and low-level intrinsic rewards results in the highest success rates and fastest convergence.

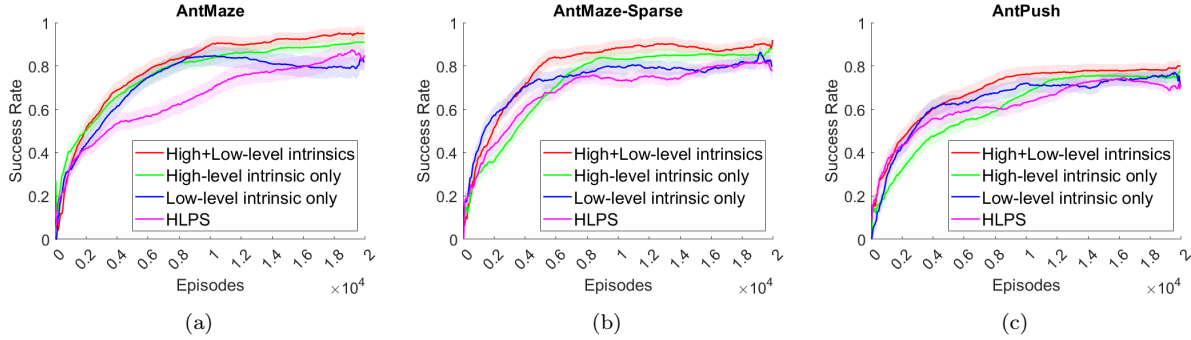


Figure 5: Success Rate on (a) AntMaze (b) AntMaze-Sparse and (c) AntPush using HLPS-G4RL, HLPS + High-level intrinsic, HLPS + Low-level intrinsic and HLPS. All curves have been equally smoothed for better visualization. The combination of high-level and low-level intrinsic rewards results in the highest success rates and fastest convergence.

on exploration, primarily refining the execution of local behaviors. These results demonstrate that intrinsic rewards at different hierarchy levels serve complementary functions, and their combination yields superior performance.

To assess the trade-off between computational efficiency and performance, we evaluate two acceleration strategies mentioned in Section 3.5. First, we vary the sampling frequency of node features by testing intervals of 1, 5, and 10 steps in HLPS. As shown in Figure 6, increasing the sampling interval substantially reduces computation time, as it decreases the number of interactions with the graph, with only minor degradation in success rates across both AntMaze and AntPush tasks. Second, we vary the proportion of training data used for the graph encoder-decoder, testing 50%, 75%, and 100% subsets. Results in Figure 7 indicate that reducing the amount of training data leads to only marginal improvements in computational efficiency and has negligible impact on final performance.

These findings suggest that the primary computational bottleneck of G4RL lies in the graph construction and node comparison processes described in Section 3.1.1, rather than in the encoder-decoder training itself. Adjusting the sampling frequency is therefore an effective approach for reducing time cost while largely preserving the benefits of G4RL integration.

### 3.7 Conclusion

This work has presented a novel approach using a graph encoder-decoder to address the challenges of poor subgoal representations and sample inefficiency in GCHRL. The proposed architecture is designed to efficiently evaluate unseen states by operating in the graph representation space. It is easy to implement and can



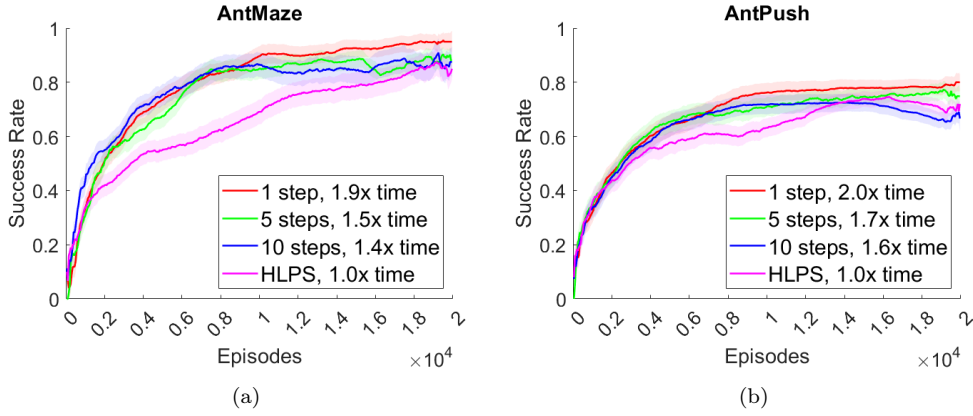


Figure 6: Success Rate on (a) AntMaze and (b) AntPush using HLPS+G4RL and HLPS. The number of steps in the legend indicates the selection of  $t_c$  as described in Section 3.5 and the timescale is calculated w.r.t. the vanilla HLPS algorithm. Increasing the sampling interval substantially reduces computation time, with only minor degradation in success rates across both AntMaze and AntPush tasks.

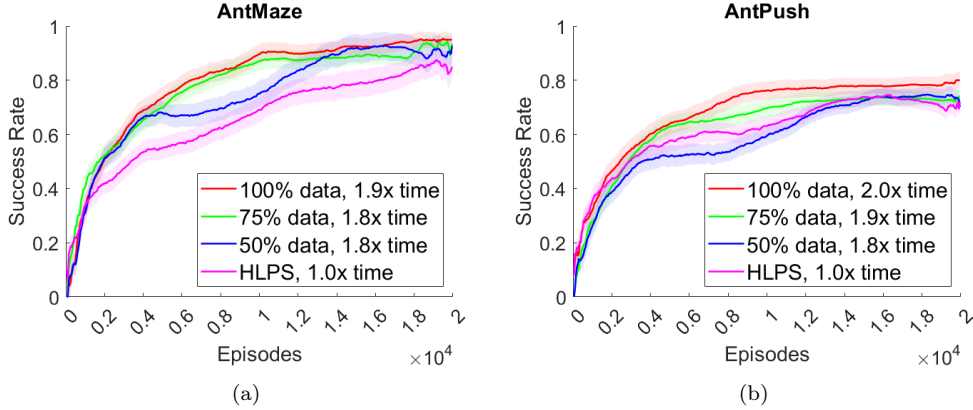


Figure 7: Success Rate on (a) AntMaze and (b) AntPush using HLPS+G4RL and HLPS. The percentage of data in the legend indicates the amount of data used in the training of the graph encoder-decoder as described in Section 3.5 and the timescale is calculated w.r.t. the vanilla HLPS algorithm. Reducing the amount of training data leads to only marginal improvements in computational efficiency and has negligible impact on final performance.

be seamlessly integrated into any existing GCHRL algorithms to enhance their performance in primarily symmetric environments. The experiments on both sparse and dense control tasks have demonstrated the effectiveness and robustness of the method.

Despite the advantages demonstrated by G4RL in GCHRL tasks, its effectiveness depends significantly on several hyperparameters (e.g.,  $\epsilon_d$ ,  $\alpha_l$ , and  $\alpha_h$ ), which require careful tuning to achieve optimal performance across different environments. In future work, we aim to develop methods for automatically selecting these hyperparameters based on environmental dynamics, thereby reducing the need for manual tuning. Also, we plan to extend our work by exploring how to generate subgoals with more interpretable representations to facilitate knowledge transfer, potentially leveraging alternative graph representations (e.g. Graph Laplacian). Another promising direction is to transfer the knowledge embedded in the graph structure to new tasks by analyzing graph topology and establishing mappings between nodes of different state graphs. Additionally, more evidence on asymmetric environments is needed to demonstrate G4RL’s robustness when asymmetric/irreversible transitions are present.

## 4 Shapley Credit Assignment for Denser RLHF

This section details our second published work, Shapley Credit Assignment Rewards (SCAR) [11]. We first review the standard RLHF setup and the inherent reward sparsity problem. We then introduce a game-theoretic perspective on text generation for credit assignment, define SCAR based on Shapley values, and finally discuss efficient approximation techniques.

### 4.1 RLHF and Reward Sparsity

We formulate the autoregressive text generation process using the standard Markov Decision Process (MDP) formalism, denoted by  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, R, \gamma)$ . The process begins in an initial state  $s_0 \in \mathcal{S}$ , which corresponds to the input prompt  $x$ . At each discrete time step  $t \in \{0, 1, \dots, T-1\}$ , the agent (the language model) is in state  $s_t$ , representing the concatenation of the prompt and the sequence generated thus far:  $s_t = x \oplus y_{1:t}$  (where  $y_{1:0}$  denotes an empty sequence, so  $s_0 = x$ , and  $\oplus$  denotes token concatenation). The agent selects an action  $a_t \in \mathcal{A}$ , which corresponds to choosing the next token  $y_{t+1}$  from the vocabulary  $\mathcal{V}$  according to its policy  $\pi_\theta(a_t|s_t)$ , parameterized by  $\theta$ . We identify the action space with the vocabulary, i.e.,  $\mathcal{A} = \mathcal{V}$ . The state transition function  $P(s_{t+1}|s_t, a_t)$  is deterministic in this context. Upon taking action  $a_t = y_{t+1}$  in state  $s_t$ , the system transitions to next state  $s_{t+1} = s_t \oplus a_t = x \oplus y_{1:t+1}$ . This generation process continues until a maximum sequence length  $T$  is reached or a designated end-of-sequence (EOS) token is generated. For notational simplicity, we often assume a fixed horizon  $T$ . The discount factor  $\gamma$  is typically set to 1 in finite-horizon text generation tasks.

In standard RLHF pipeline [14, 44, 37], a reward model  $R_\phi$ , parameterized by  $\phi$ , is trained beforehand on a dataset of human preferences  $\mathcal{D}_{\text{pref}} = \{(y^w, y^l)_i\}$ , where  $y^w$  is preferred over  $y^l$  by human annotators. The reward model assigns a scalar score  $r_\phi(x, y)$  reflecting the quality or preference level of a *completed* sequence  $y$  given the prompt  $x$ .

To stabilize training and prevent the policy  $\pi_\theta$  from drifting too far from a reference distribution (often the initial pre-trained LLM, denoted  $\pi_{\text{ref}}$ ), the reward signal used for optimization is typically augmented with a penalty term at each step  $t$ . A common choice is the Kullback-Leibler (KL) divergence between the current policy  $\pi_\theta$  and the reference policy  $\pi_{\text{ref}}$ . The standard objective is often formulated as:

$$\mathcal{J}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta} \left[ \sum_{t=1}^T R_t^{\text{orig}} \right] \quad (11)$$

where  $\mathcal{D}$  is the dataset and the reward at each timestep  $t$  in the standard sparse setup is given by

$$R_t^{\text{orig}} = R_t^{\text{KL}} + \mathbb{I}(t = T) \cdot r_\phi(x, y) \quad (12)$$

Here,  $R_t^{\text{KL}} = -\beta \log(\pi_\theta(y_t|x, y_{<t})/\pi_{\text{ref}}(y_t|x, y_{<t}))$  is the KL penalty associated with timestep  $t$ ,  $\beta$  is the KL coefficient, and  $\mathbb{I}(t = T)$  is an indicator function ensuring the reward model score  $r_\phi(x, y)$  is assigned only at the final timestep  $T$ . This terminal assignment makes the reward signal inherently *sparse*, posing a significant challenge for credit assignment during RL training. Such sparsity directly undermines the efficiency of the learning process and frequently leads to the convergence to suboptimal policies [34, 8, 49]. To overcome this, a principled approach to redistribute the terminal reward more densely across the generation steps is desirable.

### 4.2 A Game-Theoretic Framework for Credit Assignment

We frame the generation of a sequence  $y$  for a given prompt  $x$  as a cooperative game. Let the generated text  $y$  be segmented into  $N$  contiguous units,  $y = (u_1, u_2, \dots, u_N)$ . These units could be tokens, spans, sentences, or paragraphs depending on the task. The “players” in this game are these  $N$  text units. Let  $\mathcal{P} = \{u_1, \dots, u_N\}$  be the set of players.

**Characteristic Function.** The value of cooperation among a subset (coalition)  $S \subseteq \mathcal{P}$  of players is defined by a characteristic function  $v : 2^{\mathcal{P}} \rightarrow \mathbb{R}$ . This function should quantify the collective contribution of the units in  $S$  towards the final reward objective. We define  $v(S)$  based on the reward model’s evaluation of the partial text sequence formed by concatenating the units  $\{u_i \mid u_i \in S\}$  in their original order. Let  $y_S$  denote this concatenated partial sequence. Then, the value function is defined as:

$$v(S) = r_\phi(x, y_S) \quad (13)$$

For the empty set,  $v(\emptyset) = 0$ . The value of the grand coalition  $v(\mathcal{P})$  corresponds to the original sparse reward for the complete sequence,  $v(\mathcal{P}) = r_\phi(x, y)$ . Note that evaluating  $r_\phi$  on partial sequences  $y_S$  requires careful consideration, as  $y_S$  represents an incomplete sequence. Ideally,  $v(S)$  could represent the expected reward obtained by keeping the units in  $S$  fixed and sampling the remaining units from the current policy  $\pi_\theta$ . In the implementation, we resort to a practical approximation. To evaluate  $v(S)$ , we construct a sequence in which the tokens belonging to units  $u_i \in S$  are placed in their original order. The positions corresponding to units  $u_j \notin S$  are filled using empty spaces.

**Shapley Value Calculation.** The Shapley Value  $SV_{u_i}(v)$  for a player  $u_i \in \mathcal{P}$  (text unit  $u_i$ ) quantifies its fair contribution to the grand coalition value  $v(\mathcal{P})$ , calculated as the average marginal contribution of player  $u_i$  over all possible permutations of player arrivals:

$$SV_{u_i}(v) = \sum_{S \subseteq \mathcal{P} \setminus \{u_i\}} \frac{|S|!(N - |S| - 1)!}{N!} [v(S \cup \{u_i\}) - v(S)]. \quad (14)$$

The Shapley values uniquely satisfies axioms such as *efficiency* ( $\sum_{u_i \in \mathcal{P}} SV_{u_i}(v) = v(\mathcal{P})$ ), *symmetry* (equal reward for equal contribution), *linearity*, and the *null player* property (no contribution means no credit), making it a principled choice for fair credit allocation [43].

### 4.3 Shapley Values as Dense Rewards

We use the calculated Shapley values to define a dense reward signal for the RL agent. Let unit  $u_i$  consist of tokens generated from timestep  $t_{i-1} + 1$  up to and including timestep  $t_i$  (with  $t_0 = 0$ , so  $t_i$  marks the completion timestep of unit  $u_i$ ). We assign the Shapley value  $SV_{u_i}(v)$  associated with unit  $u_i$  as an additional reward component specifically at the timestep  $t_i$  marking the completion of that unit. Let  $R_t^{\text{shap}}$  denote this Shapley-based reward at timestep  $t$ . Then,

$$R_t^{\text{shap}} = \begin{cases} SV_{u_i}(v) & \text{if } t = t_i \text{ for some unit } u_i \\ 0 & \text{otherwise} \end{cases} \quad (15)$$

This component distributes the total reward  $r_\phi(x, y)$  across the episode based on the Shapley contributions, since  $\sum_{t=1}^T R_t^{\text{shap}} = \sum_{i=1}^N SV_{u_i}(v) = r_\phi(x, y)$  (due to the efficiency property, where  $N$  is the total number of units).

We then define the total reward  $R_t$  provided to the RL agent at timestep  $t$  as a convex combination of the dense Shapley-based credit allocation and the original sparse terminal reward, while retaining the per-step KL penalty. Using a hyperparameter  $\alpha \in [0, 1]$ , the total reward is:

$$R_t(\alpha) = R_t^{\text{KL}} + \alpha \cdot R_t^{\text{shap}} + (1 - \alpha) \cdot \mathbb{I}(t = T) \cdot r_\phi(x, y) \quad (16)$$

Here,  $\alpha$  controls the interpolation:

- If  $\alpha = 0$ ,  $R_t(0) = R_t^{\text{orig}}$ , recovering the standard sparse reward signal.
- If  $\alpha = 1$ ,  $R_t(1) = R_t^{\text{KL}} + R_t^{\text{shap}}$ . The terminal reward  $r_\phi(x, y)$  is entirely replaced by an equivalent value distributed densely according to Shapley contributions throughout the episode.
- If  $0 < \alpha < 1$ , the agent receives both the intermediate Shapley-based rewards and a residual portion of the original terminal reward.

This formulation allows flexible control over the density of the reward signal, balancing immediate feedback with the final outcome signal.

**Theorem 1** (Policy Invariance under SCAR Reward Shaping). *Consider a parameterized language model  $\pi_\theta$  with a learned reward model  $R_\phi$ . Let  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, R_t^{\text{orig}}, \gamma)$  be the original MDP with its reward from the reward model and  $\widehat{\mathcal{M}} = (\mathcal{S}, \mathcal{A}, P, R_t(\alpha), \gamma)$  be the MDP with dense Shapley reward. If  $\pi_\theta$  is optimal for  $\widehat{\mathcal{M}}$ , then  $\pi_\theta$  is also optimal for  $\mathcal{M}$ , and vice versa.*

## 4.4 Efficient Approximation of Shapley Values

The direct calculation of Shapley values using Equation equation 14 necessitates evaluating the characteristic function  $v(S)$  (defined in Eq. equation 13) for all  $2^N$  possible coalitions  $S$  of the  $N$  text units. This exponential complexity renders exact computation practically infeasible for typical sequence lengths encountered in text generation [43]. To make SCAR practical, we employ two key strategies: adaptive segmentation of text into units and efficient approximation of their Shapley values.

**Adaptive Text Segmentation as Players.** The definition of “players” (text units  $u_i$ ) in the cooperative game is crucial for both interpretability and computational tractability. We adapt the granularity of these units based on the task and the typical length of the generated responses, aiming to keep the number of players  $N$  manageable. We experiment with three levels of segmentation:

- **Token-level:** For tasks producing very short responses, each token  $y_t$  can be treated as an individual player  $u_i$ . This offers the finest granularity but results in a larger  $N$ .
- **Span-level:** For medium-length responses, we leverage constituency parsing [32] to establish a hierarchical grammatical structure over the generated sequence  $y$ . This process yields a constituency tree where tokens (leaf nodes) are organized into hierarchically nested constituents. Players are then defined as these syntactic constituents (e.g., noun phrases, verb phrases), formed by grouping tokens that share a common parent or ancestor node within this tree. This approach reduces  $N$  while preserving semantic coherence within each player unit, as constituents are inherently meaningful linguistic units.<sup>2</sup>
- **Sentence-level:** For tasks generating longer, multi-sentence responses, each sentence in the output  $y$  constitutes a player. Segmentation is achieved using standard sentence boundary detection. This approach markedly reduces  $N$ , especially for verbose outputs.

The choice of segmentation strategy is a hyperparameter, allowing a trade-off between the granularity of credit assignment and the computational cost of Shapley value estimation.

**Approximation Using Owen Values.** To ensure the practical applicability of SCAR, we employ an approximation scheme based on Owen value [38], which is a coalitional extension of Shapley values designed for games where players are grouped into a predefined coalition structure[4]. For the task, a hierarchical structure  $\mathcal{B}$  is imposed on the sequence of  $N$  text units, achieved by applying a heuristic parsing algorithm to the units. This partition  $\mathcal{B}$  defines nested groupings of the units. The Owen value is then computed with respect to this partition  $\mathcal{B}$ .

Let  $N$  be the set of all players, and let  $v : 2^N \rightarrow \mathbb{R}$  be the characteristic function. The **coalition structure** is a partition of  $N$ , denoted by  $\mathcal{B} = \{B_1, B_2, \dots, B_m\}$ , where  $m = |\mathcal{B}|$  is the number of unions. For any player  $i \in N$ , let  $B(i)$  be the unique union in  $\mathcal{B}$  that contains player  $i$ . The Owen value  $\Phi_i(v, \mathcal{B})$  for player  $i$  is given by the formula:

$$\Phi_i(v, \mathcal{B}) = \sum_{R \subseteq \mathcal{B} \setminus \{B(i)\}} \frac{|R|!(m - |R| - 1)!}{m!} \cdot \text{MC} \quad (17)$$

Where MC is the marginal contribution, defined as:

$$\text{MC} = v \left( \bigcup_{B \in R} B \cup \{B(i)\} \right) - v \left( \bigcup_{B \in R} B \right) \quad (18)$$

Marginal contributions are evaluated by forming coalitions structurally: combinations involving subsets within a unit’s own group are explored, while units belonging to other groups in the partition are treated as indivisible blocks, as they are either entirely included or entirely excluded from a coalition, rather than exploring all their individual subsets. By limiting the evaluation to coalitions dictated by the partition structure  $\mathcal{B}$ , the number of required characteristic function evaluations (reward model queries) is substantially reduced compared to the exact Shapley computation. Consequently, the computational complexity is reduced from exponential,  $O(2^N)$ , to quadratic in  $N$ , rendering the approach tractable. We use the SHAP package [29] for Shapley values and Owen values computation.

<sup>2</sup><https://www.nltk.org/howto/parse.html>

## 4.5 Results

In this section, we empirically evaluate the effectiveness of SCAR across three distinct tasks characterized by varying response lengths. My primary objective is to demonstrate that SCAR enables more efficient and effective training compared to standard sparse RLHF and alternative dense reward baselines.

<b>Prompt:</b>	“While some scenes were”												
<b>Sparse:</b>	initially	disturbing	to	sit	through,	they	ultimately	contributed	to	a	deeply	powerful	and
	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	moving story . <EOS>												
	0.0	0.0	0.0	+10.0									
<b>ABC:</b>	initially	disturbing	to	sit	through,	they	ultimately	contributed	to	a	deeply	powerful	and
	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	+0.8	+0.8	+0.2
	moving story . <EOS>												
	+1.3	+1.6	+5.3	0.0									
<b>SCAR:</b>	initially	disturbing	to	sit	through,	they	ultimately	contributed	to	a	deeply	powerful	
	+0.8	-0.5	-0.1	-1.8	-0.8	-0.2	+1.0	+0.7	0.0	+0.5	+1.8	+2.2	
	and moving story . <EOS>												
	+0.5	+4.0	+0.7	+0.8	0.0								

Figure 8: Comparison of reward distribution strategies for an example generated sequence. Sparse RLHF assigns the total reward at the end. SCAR and ABC distribute this reward across tokens/spans based on their respective methodologies, shown with background highlights (color hue for sign, intensity for magnitude; more intense/saturated means higher absolute contribution) and numerical scores.

**Evaluation Tasks and Models.** We evaluate the proposed method across three diverse tasks prevalent in RLHF research: sentiment control, text summarization, and instruction tuning [13, 50]. For sentiment control and instruction tuning, we utilize the implementation in [13]. However, for summarization, due to difficulties in reproducing the results, we switched to the implementation in [22]. The datasets, policy models, and reward models used for each task are described in more detail below.

**Sentiment Control:** The objective is to generate positive reviews of movies. We use the IMDB dataset [31]. The policy model is GPT-2 small [40], initialized by fine-tuning for one epoch on the IMDB training set. During RLHF training, prompts are generated by randomly selecting the first 4 to 8 tokens from reviews in the training data. The reward signals are provided using a pre-trained sentiment classifier, same as [13].

**Text Summarization:** We evaluate the method on the automatic text summarization task, following prior work [44, 13, 26]. For this evaluation, we use the Reddit TL;DR dataset [48], specifically the filtered version [44], which includes approximately 116K training examples, 6K validation examples, and 6K test examples. The policy model used is Pythia-1B [9], which we initialize via supervised fine-tuning on the training set for 2,000 steps with a batch size of 64. Additionally, we train a 1B-parameter reward model (initialized using the SFT model) on 92K human preference pairs, achieving approximately 74% accuracy on the validation set.

**Instruction Turning:** We evaluate models on the task of following user instructions. To do this, we fine-tune language models using the helpfulness subset of the Anthropic Helpful and Harmless (HH) dataset [7], which contains 43K human-written prompts paired with model responses that have been ranked by human annotators for helpfulness. Preference is based on which response is more informative and helpful for the task. The policy model is initialized using the OpenLLaMA-7B model [20], an open-source reproduction of Meta’s LLaMA collection [47] trained on fully open-source dataset. For the reward model, we use a 3B reward model [17]. This reward model was trained using the same HH dataset, where it learns to assign a scores to candidate completions based on their predicted usefulness.

**Baselines.** We compare the proposed method against three key baselines. **RLHF** represents the standard approach using the sparse terminal reward (Eq. 12) with KL regularization. **ABC (Attention Based Credit)** [13] uses reward model attention scores for dense rewards distribution. **Uniform** is a baseline that distributes the terminal reward evenly across all tokens. For fair comparison, all methods are optimized using

the PPO objective [42] with consistent hyperparameters. All methods initialize their policy models using the same SFT checkpoint to ensure a common starting point. All experiments were conducted on a single A100 GPU (80GB VRAM), and results are averaged over 5 random seeds.

**Evaluation Metrics.** We track the average reward  $r_\phi(x, y)$  per episode during training to evaluate learning speed and the level of convergence. The final performance is reported as the mean reward on the test set after convergence. For the summarization task, we additionally employ **LLM-as-a-judge** [53] evaluation to compare the quality (e.g., accuracy, coverage, conciseness, clarity, and coherence) of summaries generated by models trained with different methods. We randomly sampled 1K summaries from the TL;DR test set for LLM evaluation. To mitigate potential positional bias in these pairwise comparisons, we randomize the presentation order of summaries. For the instruction-tuning task, we use AlpacaEval [28] to compare the quality of 1K model’s response. AlpacaEval is designed to better handle potential issues such as length bias, thereby providing a more reliable assessment of response quality.

Figure 8 provides a qualitative illustration of how SCAR distributes rewards compared to sparse RLHF and ABC for an example generated sequence. Sparse RLHF, by definition, assigns the entire reward only at the end of the sequence. ABC, which uses attention scores from the reward model’s final layer to distribute rewards, tends to concentrate rewards on tokens near the end of the sequence. As seen in the example, significant credit is assigned to the final punctuation mark (“.”), while earlier, crucial tokens receive almost zero attention scores. Furthermore, standard attention scores are non-negative, making it difficult for ABC to assign explicit negative credit to tokens or spans that detract from the output quality. For instance, a phrase like “*disturbing to sit through*” that negatively impact the perceived sentiment, would not receive negative rewards from ABC. In contrast, SCAR can assign both positive and negative rewards to tokens based on their game-theoretic marginal contribution.

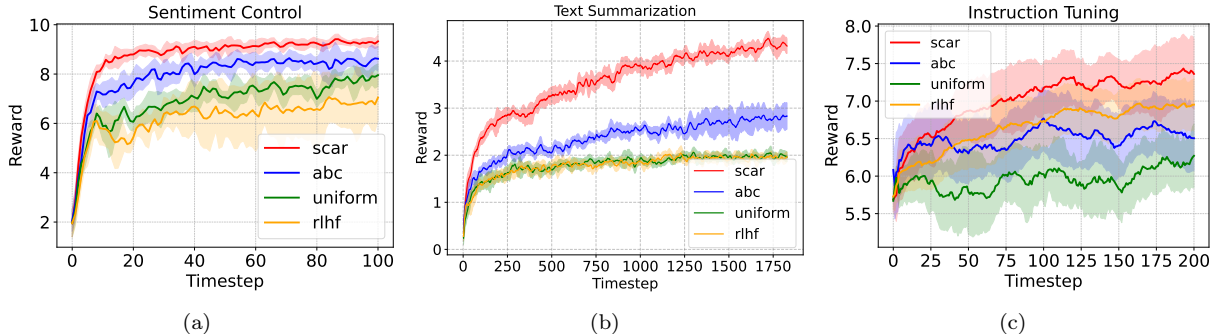


Figure 9: Average reward per timestep during RLHF training for sentiment control (left), text summarization (center), and instruction tuning (right). Curves show the mean reward across five random seeds, with shaded regions representing the standard deviation. SCAR consistently demonstrates faster convergence and achieves higher or comparable final reward levels compared to sparse RLHF, Uniform reward distribution, and Attention-Based Credit (ABC) baselines.

Task	Sparse RLHF	Uniform	ABC	SCAR
IMDB	6.86 ± 0.86	7.73 ± 0.02	8.48 ± 1.60	<b>9.27 ± 0.00</b>
TL;DR	1.60 ± 0.11	1.68 ± 0.02	2.85 ± 0.21	<b>4.35 ± 0.11</b>
HH-RLHF	6.93 ± 0.00	6.17 ± 0.00	6.59 ± 0.01	<b>7.31 ± 0.01</b>

Table 1: Average reward scores for the trained policy on the test sets for sentiment control (IMDB), text summarization (TL;DR), and instruction tuning (Anthropic HH). Higher scores indicate better performance. Results are averaged over 5 random seeds. Best performance per task is in **bold**.

As depicted in Figure 9, SCAR consistently demonstrated advantages over three baseline methods in terms of learning speed and convergence across all three tasks. This consistent pattern across diverse tasks suggests that the principled, Shapley value-based credit assignment offered by the method effectively improves learning efficiency and enhances the final policy performance. Table 1 presents the average reward scores achieved by each method on the held-out test sets for the three tasks. As shown in the table, SCAR-tuned policy consistently achieves the highest performance compared to the baselines.

For summarization, we use **gemini-2.5-pro** to compare anonymized model outputs (SCAR vs. baselines) on coherence, relevance, conciseness, and overall quality. As shown in Table 2, summaries generated by the SCAR-tuned model were preferred over those from the ABC-tuned model in 60.3% and over the sparse RLHFC-tuned model in 61.2%. For the instruction tuning task, we leveraged AlpacaEval [28] and used **gpt-4-turbo** as an LLM judge, to conduct robust pairwise evaluations. These evaluations assessed helpfulness, harmlessness, and adherence to instructions. SCAR-generated responses achieved a win rate of 54.9% when compared against ABC, and a win rate of 56.3% against standard sparse RLHF. These results provide further evidence that the improvements from SCAR translate to genuinely higher-quality outputs according to human-like preferences.

Baselines	Win (%)
<i>Text Summarization (Reddit TL;DR)</i>	
vs. RLHF	61.2%
vs. ABC	60.3%
<i>Instruction Tuning (Anthropic HH)</i>	
vs. RLHF	56.3%
vs. ABC	54.9%

Table 2: LLM-as-Judge pairwise win rates for SCAR against baselines.

## 4.6 Conclusion

This part has presented a novel approach using a graph encoder-decoder to address the challenges of poor subgoal representations and sample inefficiency in GCHRL. The proposed architecture is designed to efficiently evaluate unseen states by operating in the graph representation space. It is easy to implement and can be seamlessly integrated into any existing GCHRL algorithms to enhance their performance in primarily symmetric environments. The experiments on both sparse and dense control tasks have demonstrated the effectiveness and robustness of the method.

Despite its strengths, SCAR has limitations: the computational overhead of Shapley approximations, even with optimizations like Owen values and adaptive segmentation; the assumption that the reward model can meaningfully score partial sequences, which may not suit certain types like rule-based models that only evaluate final answers (e.g., in mathematical reasoning). Future work will target more efficient approximation techniques, robust and adaptive segmentation methods, and rigorous evaluation on larger-scale language models and broader tasks.

## 5 Proposed Research: Dense reward from previous knowledge for continual RL

The ability to continually learn and adapt to new environments [23] while retaining knowledge of previously visited ones is crucial for modern RL algorithms. Previous methods focus on mitigating plasticity loss in neural networks [1, 15, 16]. Recently, people have tried to develop RL agents which can acquire structured knowledge and adapt to new environments [3]. This requires that the agent can 1) store the knowledge acquired in previous environments in a structured, transferable form and 2) adapt the stored knowledge to the new environment/domain easily. When the agent learns in tasks with similar dynamics, this similarity in task structure can be considered as "invariant knowledge" and can be used to transfer. To achieve this, we naturally consider graphs as the carriers of transferable knowledge, which are suitable for modelling environmental dynamics. The approximated state graphs model the environmental dynamics in terms of node connectivity and degrees. Furthermore, the information in the graph nodes/labels is easy to transfer to the new environment.

**Graph construction in the first environment:** The graph construction process is similar to what we discussed in Section 3.1.1. Suppose the representation is different from any state representations stored in the graph. In that case, we store the state representation  $\phi(s_t)$  as the node feature of an empty node in the graph and build an edge between this node and the node that corresponds to the previous state. In addition,

we store the estimated current state value into the node label:

$$\forall s_v \in \mathcal{V}, \|\phi(s_t) - \phi(s_v)\|_2 > \epsilon_d, \quad (19)$$

$$\mathbf{A}_{\phi(s_t), \phi(s_{t-1})} = \mathbf{A}_{\phi(s_{t-1}), \phi(s_t)} = 1, \quad \mathbf{y}_{s_t} = V_\pi(s_t), \quad (20)$$

where  $\mathbf{y}_{s_t}$  represents the state value of  $s_t$ .

Graph updating in the first environment: Suppose the graph is now full. When a new state  $s_t$  is encountered,

$$s_v = \arg \min_{s_u: \|\phi(s_t) - \phi(s_u)\|_2 \leq \epsilon_d} \|\phi(s_t) - \phi(s_u)\|_2. \quad (21)$$

If  $s_v$  exists, as before we relabel the node as  $\phi(s_t)$  and evaluate the state value:

$$\mathbf{y}_{s_t} = V_\pi(s_t) \quad (22)$$

Otherwise, we replace the oldest state node in the graph with the current state node, delete all edges previously linked to that node, and create an edge  $(s_{t-1}, s_t)$  with weight  $\mathbf{A}_{\phi(s_{t-1}), \phi(s_t)} = \mathbf{A}_{\phi(s_t), \phi(s_{t-1})} = 1$  and set  $\mathbf{y}_{s_t} = V_\pi(s_t)$ .

After sufficient learning in the first environment  $E_1$ , we calculate the graph laplacian of the state graph  $\mathbf{L}$

$$\mathbf{L} = \mathbf{D} - \mathbf{A} \quad (23)$$

where  $\mathbf{D}$  is the diagonal degree matrix. For later uses, this  $\mathbf{L}$  will be denoted by  $\mathbf{L}_1$ .

**Comparing with knowledge from the first environment in the second environment** We introduce a second environment  $E_2$  where the task is designed structurally analogous to the first, yet varies in its specifics, requiring retraining within the traditional RL framework. For this environment, we construct and update a graph using the identical procedure. At regular intervals, we find the graph Laplacian  $\mathbf{L}_2$ . While canonical labelling [5] is the standard tool for verifying exact graph isomorphism, it is ineffective when graphs are only approximately isomorphic. Therefore, we employ a spectral method to quantitatively assess the degree of similarity by measuring the difference between the associated Laplacian matrices  $\mathbf{L}_1$  and  $\mathbf{L}_2$ . We compute the eigen decompositions:

$$\mathbf{L}_k = \mathbf{V}_k \mathbf{\Lambda}_k \mathbf{V}_k^T = \lambda_k^{(1)} \mathbf{v}_k^{(1)} (\mathbf{v}_k^{(1)})^T + \lambda_k^{(2)} \mathbf{v}_k^{(2)} (\mathbf{v}_k^{(2)})^T + \dots + \lambda_k^{(N)} \mathbf{v}_k^{(N)} (\mathbf{v}_k^{(N)})^T, \quad k = 1, 2, \quad (24)$$

where  $\mathbf{V}_k = [\mathbf{v}_k^{(1)}, \mathbf{v}_k^{(2)}, \dots, \mathbf{v}_k^{(N)}]$  is orthogonal and the largest entry of each eigenvector  $\mathbf{v}_k^{(i)}$  in magnitude is positive, and  $\mathbf{\Lambda}_k = \text{diag}(\lambda_k^{(i)})$  with  $\lambda_k^{(1)} \geq \lambda_k^{(2)} \geq \dots \geq \lambda_k^{(N)} = 0$ . Each row of  $\mathbf{V}_k$  is a feature vector of the corresponding node of the graph.

Note that the two graphs are isomorphic if and only if there exists a permutation matrix  $\mathbf{P}$  such that  $\mathbf{L}_2 = \mathbf{P} \mathbf{L}_1 \mathbf{P}^T$ . To check if the two graphs are isomorphic, we first check if the difference between the two spectra is within a preset threshold:

$$\sum_{j=1}^N |\lambda_1^{(j)} - \lambda_2^{(j)}|^2 \leq \epsilon_\lambda \quad (25)$$

If it holds, we try to match feature vectors of nodes in the two graphs. There are two cases:

1. All eigenvalues are distinct. In this case, the eigenvectors  $\mathbf{v}_k^{(i)}$  are unique up to the sign. We try to match rows of  $\mathbf{V}_1$  with rows of  $\mathbf{V}_2$ . Specifically, for the first row vector of  $\mathbf{V}_1$ , we check if there exists a row vector of  $\mathbf{V}_2$  that is close to it:

$$\exists \mathbf{V}_2(i, :) \text{ such that } \|\mathbf{V}_1(1, :) - \mathbf{V}_2(i, :)\|_2 \leq \epsilon_v, \quad (26)$$

where  $\epsilon_v$  is a predefined tolerance threshold. If the match succeeds, we continue to try to match other rows of  $\mathbf{V}_1$  and  $\mathbf{V}_2$ . If one match fails, we conclude that the two graphs are not similar.

2. Some eigenvalues are repeated. This case is much more complicated. However, there exist  $O(N^3)$  practical methods to do the match, see, e.g., [18, 25]



Once two graphs are considered similar, when we act  $a_t$  in environment 2 and transition from  $s_t$  to  $s_{t+1}$  happens, we find the closest state feature to  $\phi(s_{t+1})$  stored in graph 2, denoted by  $\phi(s_{e2})$  which is the state representation of  $s_{e2}$ , then find its pairing node in graph 1 (denoted by  $s_{e1}$ ). The intrinsic reward of this transition is then defined as:

$$R_{int}(s_t, a_t, s_{t+1}) = \beta \mathbf{y}_{s_{e1}} = \beta V_{\pi}(s_{e1}) \quad (27)$$

where  $\beta$  is a hyperparameter controlling the significance of the intrinsic term.

**Ongoing Work** To complete this research project, several key areas require further investigation and development:

- **Identifying appropriate test environments:** The proposed method is designed to identify and leverage similarities between tasks. To effectively evaluate this capability, it is crucial to select or construct environments that exhibit meaningful task similarities alongside non-trivial differences. For example, we are currently experimenting with variants of the same environment that differ in granularity, or that share the same state space while varying the action space. Careful environment design will be essential to demonstrate the generality and robustness of our approach.
- **Establishing relevant baselines:** A rigorous comparison with existing continual reinforcement learning (CRL) methods is necessary to validate our approach. We plan to benchmark against recent baselines such as those proposed by [3] and [1], using our curated sets of related tasks. Selecting strong and diverse baselines will help clarify where our method stands in terms of performance and generalization.
- **Defining evaluation metrics:** In addition to standard CRL performance metrics (e.g., average return or task score), it may be useful to develop task similarity metrics or diagnostics tailored to our method. These could provide deeper insights into when and why our approach succeeds or fails, and guide future improvements.

## References

- [1] Zaheer Abbas, Rosie Zhao, Joseph Modayil, Adam White, and Marlos C Machado. Loss of plasticity in continual deep reinforcement learning. In *Conference on lifelong learning agents*, pages 620–636. PMLR, 2023.
- [2] Amr Ahmed, Nino Shervashidze, Shravan Narayanamurthy, Vanja Josifovski, and Alexander J Smola. Distributed large-scale natural graph factorization. In *Proceedings of the 22nd international conference on World Wide Web*, pages 37–48, 2013.
- [3] Nishanth Anand and Doina Precup. Prediction and control in continual reinforcement learning. *Advances in Neural Information Processing Systems*, 36:63779–63817, 2023.
- [4] Robert J Aumann and Jacques H Dreze. Cooperative games with coalition structures. *International Journal of game theory*, 3:217–237, 1974.
- [5] László Babai and Eugene M Luks. Canonical labeling of graphs. In *Proceedings of the fifteenth annual ACM symposium on Theory of computing*, pages 171–183, 1983.
- [6] Dzmitry Bahdanau, Philemon Brakel, Kelvin Xu, Anirudh Goyal, Ryan Lowe, Joelle Pineau, Aaron Courville, and Yoshua Bengio. An actor-critic algorithm for sequence prediction. In *International Conference on Learning Representations*, 2017.
- [7] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- [8] Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos. Unifying count-based exploration and intrinsic motivation. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.

- [9] Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pages 2397–2430. PMLR, 2023.
- [10] Christian Buck, Jannis Bulian, Massimiliano Ciaramita, Wojciech Gajewski, Andrea Gesmundo, Neil Houlsby, and Wei Wang. Ask the right questions: Active question reformulation with reinforcement learning. In *International Conference on Learning Representations*, 2018.
- [11] Meng Cao, Shuyuan Zhang, Xiao-Wen Chang, and Doina Precup. Scar: Shapley credit assignment for more efficient rlhf. *arXiv preprint arXiv:2505.20417*, 2025.
- [12] Shaosheng Cao, Wei Lu, and Qiongkai Xu. Grarep: Learning graph representations with global structural information. In *Proceedings of the 24th ACM international on conference on information and knowledge management*, pages 891–900, 2015.
- [13] Alex James Chan, Hao Sun, Samuel Holt, and Mihaela van der Schaar. Dense reward for free in reinforcement learning from human feedback. In *Forty-first International Conference on Machine Learning*, 2024.
- [14] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [15] Shibhansh Dohare, J Fernando Hernandez-Garcia, Qingfeng Lan, Parash Rahman, A Rupam Mahmood, and Richard S Sutton. Loss of plasticity in deep continual learning. *Nature*, 632(8026):768–774, 2024.
- [16] Shibhansh Dohare, J Fernando Hernandez-Garcia, Parash Rahman, A Rupam Mahmood, and Richard S Sutton. Maintaining plasticity in deep continual learning. *arXiv preprint arXiv:2306.13812*, 2023.
- [17] Hanze Dong, Wei Xiong, Deepanshu Goyal, Yihan Zhang, Winnie Chow, Rui Pan, Shizhe Diao, Jipeng Zhang, Kashun Shum, and Tong Zhang. Raft: Reward ranked finetuning for generative foundation model alignment. *arXiv preprint arXiv:2304.06767*, 2023.
- [18] Zhou Fan, Cheng Mao, Yihong Wu, and Jiaming Xu. Spectral graph matching and regularized quadratic relaxations: Algorithm and theory. In *International conference on machine learning*, pages 2985–2995. PMLR, 2020.
- [19] Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *International conference on machine learning*, pages 1587–1596. PMLR, 2018.
- [20] Xinyang Geng and Hao Liu. Openllama: An open reproduction of llama, May 2023.
- [21] Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*, 2018.
- [22] Shengyi Huang, Michael Noukhovitch, Arian Hosseini, Kashif Rasul, Weixun Wang, and Lewis Tunstall. The n+ implementation details of RLHF with PPO: A case study on TL;DR summarization. In *First Conference on Language Modeling*, 2024.
- [23] Khimya Khetarpal, Matthew Riemer, Irina Rish, and Doina Precup. Towards continual reinforcement learning: A review and perspectives. *Journal of Artificial Intelligence Research*, 75:1401–1476, 2022.
- [24] Junsu Kim, Younggyo Seo, and Jinwoo Shin. Landmark-guided subgoal generation in hierarchical reinforcement learning. *Advances in neural information processing systems*, 34:28336–28349, 2021.
- [25] Stefan Klus and Tuhin Sahai. A spectral assignment approach for the graph isomorphism problem. *Information and Inference: A Journal of the IMA*, 7(4):689–706, 2018.

- [26] Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Ren Lu, Thomas Mesnard, Johan Ferret, Colton Bishop, Ethan Hall, Victor Carbune, and Abhinav Rastogi. Rlaif: Scaling reinforcement learning from human feedback with ai feedback. 2023.
- [27] Jiwei Li, Will Monroe, Alan Ritter, Dan Jurafsky, Michel Galley, and Jianfeng Gao. Deep reinforcement learning for dialogue generation. In Jian Su, Kevin Duh, and Xavier Carreras, editors, *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1192–1202, Austin, Texas, November 2016. Association for Computational Linguistics.
- [28] Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. AlpacaEval: An automatic evaluator of instruction-following models. [https://github.com/tatsu-lab/alpaca\\_eval](https://github.com/tatsu-lab/alpaca_eval), 5 2023.
- [29] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 4765–4774. Curran Associates, Inc., 2017.
- [30] Yu Luo, Tianying Ji, Fuchun Sun, Huaping Liu, Jianwei Zhang, Mingxuan Jing, and Wenbing Huang. Goal-conditioned hierarchical reinforcement learning with high-level model approximation. *IEEE Transactions on Neural Networks and Learning Systems*, 2024.
- [31] Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In Dekang Lin, Yuji Matsumoto, and Rada Mihalcea, editors, *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA, June 2011. Association for Computational Linguistics.
- [32] Mitchell P. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. Building a large annotated corpus of English: The Penn Treebank. *Computational Linguistics*, 19(2):313–330, 1993.
- [33] Ofir Nachum, Shixiang Shane Gu, Honglak Lee, and Sergey Levine. Data-efficient hierarchical reinforcement learning. *Advances in neural information processing systems*, 31, 2018.
- [34] Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *ICML*, volume 99, pages 278–287. Citeseer, 1999.
- [35] Mohammad Norouzi, Samy Bengio, Zhifeng Chen, Navdeep Jaitly, Mike Schuster, Yonghui Wu, and Dale Schuurmans. Reward augmented maximum likelihood for neural structured prediction. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.
- [36] Mingdong Ou, Peng Cui, Jian Pei, Ziwei Zhang, and Wenwu Zhu. Asymmetric transitivity preserving graph embedding. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1105–1114, 2016.
- [37] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744. Curran Associates, Inc., 2022.
- [38] Guillermo Owen. Values of games with a priori unions. In *Mathematical economics and game theory: Essays in honor of Oskar Morgenstern*, pages 76–88. Springer, 1977.
- [39] Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- [40] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019.

- [41] Seonggi Ryang and Takeshi Abekawa. Framework of automatic text summarization using reinforcement learning. In Jun'ichi Tsujii, James Henderson, and Marius Pasca, editors, *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 256–265, Jeju Island, Korea, July 2012. Association for Computational Linguistics.
- [42] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [43] Lloyd S Shapley. Stochastic games. *Proceedings of the national academy of sciences*, 39(10):1095–1100, 1953.
- [44] Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 3008–3021. Curran Associates, Inc., 2020.
- [45] Richard Stuart Sutton. *Temporal credit assignment in reinforcement learning*. University of Massachusetts Amherst, 1984.
- [46] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ international conference on intelligent robots and systems*, pages 5026–5033. IEEE, 2012.
- [47] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [48] Michael V"olske, Martin Potthast, Shahbaz Syed, and Benno Stein. TL;DR: Mining Reddit to learn automatic summarization. In *Proceedings of the Workshop on New Frontiers in Summarization*, pages 59–63, Copenhagen, Denmark, September 2017. Association for Computational Linguistics.
- [49] Zeqiu Wu, Yushi Hu, Weijia Shi, Nouha Dziri, Alane Suhr, Prithviraj Ammanabrolu, Noah A Smith, Mari Ostendorf, and Hannaneh Hajishirzi. Fine-grained human feedback gives better rewards for language model training. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 59008–59033. Curran Associates, Inc., 2023.
- [50] Eunseop Yoon, Hee Suk Yoon, SooHwan Eom, Gunsoo Han, Daniel Nam, Daejin Jo, Kyoung-Woon On, Mark Hasegawa-Johnson, Sungwoong Kim, and Chang Yoo. TLCR: Token-level continuous reward for fine-grained reinforcement learning from human feedback. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association for Computational Linguistics: ACL 2024*, pages 14969–14981, Bangkok, Thailand, August 2024. Association for Computational Linguistics.
- [51] Shuyuan Zhang, Zihan Wang, Xiao-Wen Chang, and Doina Precup. Incorporating spatial information into goal-conditioned hierarchical reinforcement learning via graph representations. *Transactions on Machine Learning Research*, 2025.
- [52] Tianren Zhang, Shangqi Guo, Tian Tan, Xiaolin Hu, and Feng Chen. Adjacency constraint for efficient hierarchical reinforcement learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4):4152–4166, 2022.
- [53] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging LLM-as-a-judge with MT-bench and chatbot arena. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023.