

Option Discovery Using LLM-guided Semantic Hierarchical Reinforcement Learning

Chak Lam Shek¹ and Pratap Tokekar²

Abstract—Large Language Models (LLMs) have shown remarkable promise in reasoning and decision-making, yet their integration with Reinforcement Learning (RL) for complex robotic tasks remains underexplored. In this paper, we propose an LLM-guided hierarchical RL framework, termed LDSC, that leverages LLM-driven subgoal selection and option reuse to enhance sample efficiency, generalization, and multi-task adaptability. Traditional RL methods often suffer from inefficient exploration and high computational cost. Hierarchical RL helps with these challenges, but existing methods often fail to reuse options effectively when faced with new tasks. To address these limitations, we introduce a three-stage framework that uses LLMs for subgoal generation given natural language description of the task, a reusable option learning and selection method, and an action-level policy, enabling more effective decision-making across diverse tasks. By incorporating LLMs for subgoal prediction and policy guidance, our approach improves exploration efficiency and enhances learning performance. On average, LDSC outperforms the baseline by 55.9% in average reward, demonstrating its effectiveness in complex RL settings. More details and experiment videos could be found in [this link](#)¹.

I. INTRODUCTION

In many decision-making scenarios, robots are tasked with efficiently managing a wide array of tasks, each exhibiting varying levels of complexity [1], [2]. Despite significant advancements in Reinforcement Learning (RL), several challenges remain, including inefficient exploration [3], [4] and high computational cost [5], [6], particularly when robots are required to adapt to new goals or environments. These challenges stem from the need to repeatedly explore similar state spaces across different tasks, which makes the learning process both time-consuming and resource-intensive. To address this issue, some approaches (e.g., [7]), advocate for transferring pretrained policies across various environments. However, such pretrained policies often fail to generalize effectively, offering minimal benefits in completely new tasks or when prior experience is insufficient [8], [9].

Numerous studies have investigated Hierarchical Reinforcement Learning (HRL) [10], which aims to learn reusable options [11] or skills [12] from tasks and apply them to new problem settings. One prominent example is the Deep Skill Chaining (DSC) [13], which hierarchically decomposes

This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible. This is the preprint version.

¹Department of Electrical Engineering, University of Maryland, College Park, USA cshel1@umd.edu

²Department of Computer Science, University of Maryland, College Park, USA tokekar@umd.edu

¹<https://raaslab.org/projects/LDSC>

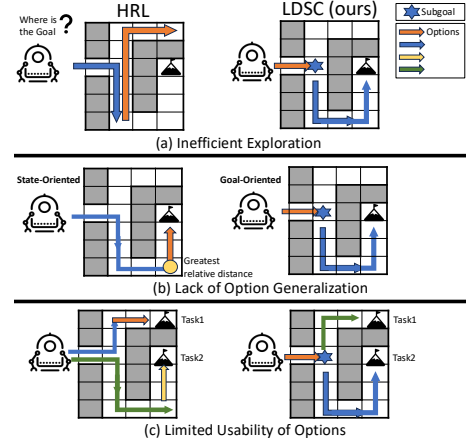


Fig. 1. Comparison of HRL and LDSC in terms of exploration, generalization, and usability. (a) Inefficient Exploration: LDSC improves exploration by leveraging subgoals to guide the agent efficiently, reducing the need for exhaustive search. (b) Lack of Option Generalization: Unlike HRL, which learns state-oriented options that do not generalize well, LDSC learns goal-oriented options that can be reused across different tasks. (c) Limited Usability of Options: LDSC structures options around subgoals, enhancing their reusability for multiple tasks, improving learning efficiency, and enabling better skill transfer.

decision-making into a policy over options² for skill selection and intra-option policies for generating low-level actions. The success of this approach lies in grouping actions into options that can be efficiently selected and executed to achieve long-term goals. However, HRL encounters three primary challenges, as illustrated in Figure 1: 1) **Inefficient Exploration** where robots face difficulties in developing effective option policies for long-horizon tasks [15] due to the need to reach goals before receiving positive feedback, leading to slow learning progress. 2) **Limited Policy Generalization** as most HRL approaches are heavily state-dependent [16], causing learned policies to be difficult to abstract or transfer to tasks with different state configurations. 3) **Restricted Usability of Learned Options** as these options are often tied to specific state transitions or distances, limiting their ability to be reused across tasks with varying goals or environments. Collectively, these challenges significantly impede HRL's capacity to generalize and adapt to new and diverse settings.

Recent advancements in decision-making have explored reasoning-based approaches, such as leveraging Large Language Models (LLMs) to generate plans [17] or provide rewards [18] that guide robot behavior. These approaches

²An *option* is defined as a tuple $(\mathcal{I}, \pi, \beta)$, consisting of an initiation set \mathcal{I} , an intra-option policy π , and a termination condition β [14].

integrate Semantic reasoning to structure decision-making, enabling robots to plan and adapt to new tasks more efficiently. By employing semantic representations, robots can reason about abstract concepts such as subgoals and actions, which is particularly beneficial when facing tasks with high complexity or unfamiliar environments.

To address the challenge of reusing option policies in multi-task settings, we propose LLM-guided Semantic Deep Skill Chaining (LDSC), a Semantic HRL method that leverages semantic logic to enable effective policy transfer across different tasks. The method follows a three-level hierarchy: the subgoal policy operates at a reasoning level, determining the sequence of subgoals required to achieve the high-level goal. This creates a logical framework that guides the option and action policies. The option policy selects the appropriate option for the chosen subgoal, while the action policy generates the corresponding low-level actions.

LDSC leverages LLM reasoning to guide the learning process and addresses several challenges in traditional HRL. By enabling the robot to reason about subgoals, LDSC provides a high-level, structured plan that alleviates the issue of inefficient exploration. Rather than relying on relative state distances, LDSC uses semantic representations that make the option policies more flexible and goal-oriented. This approach allows the robot to move meaningfully between subgoals, improving the reusability of learned options. Consequently, LDSC reduces the need for retraining and ensures that the learned options are adaptable and transferable across a variety of tasks, making it well-suited for complex, multi-task environments.

Below, we outline the three main contributions of this work:

- **Improved Learning Efficiency:** Through LLM-generated subgoals, robots can achieve structured task completion, **accelerating the learning process**.
- **Policy Generalization:** Our approach enables effective transfer of learned policies across **diverse tasks**, promoting **generalization** and **adaptability** in complex environments.
- **Experimental Validation:** We demonstrate the effectiveness of our method through extensive experiments in diverse environments, showcasing its applicability to real-world multi-task challenges. On average, **LDSC outperforms baseline methods by 55.9%**, reduces task completion time by **53.1%**, and improves success rates by **72.7%**, while maintaining a similar training time to DSC.

II. RELATED WORK

Sharma et al. [19] uses RL to discover low-level skills with predictable outcomes, allowing model-based planning in skill space. Bacon et al. [20] introduces an option-critic architecture that autonomously learns temporal abstractions (options) in RL. Chunduru et al. [21] proposes an attention-based extension that improves option diversity. Option discovery based on the construction of reusable skills through Successor Representations [22], learning from small fragments of

experience within options [23], the discovery of high-level behaviors from low-level actions [24], unsupervised skill emergence from information-theoretic objectives [25], and the use of representations encoding diffusive information flow [16]. Skill chaining [15] is a method for discovering new skills by creating chains of skills. Deep Skill Chaining [13], an extension of skill chaining enables the discovery of useful skills in high-dimensional spaces. In a related direction, [26] studies combinatorial generalization in RL and proposes temporal data augmentation to enable policy generalization across unseen state-goal pairs

Recently, numerous papers have explored the integration of LLMs for reasoning and their application in improving RL. Recent studies have demonstrated the use of LLMs for task decomposition into subproblems, including [27], [28], [29]. LLM-Planner [30], [17] leverages LLMs to generate and adapt high-level plans using natural language instructions and environmental observations. Chain-of-thought methods have been explored to improve reasoning in LLMs [31], [32]. The use of tree structures for planning has also been studied in works such as [33], [34], [35]. Several works have explored using LLMs to convert natural language instructions into reward signals, including [36], [37]. The CLIP model has been applied as a zero-shot reward model, assigning rewards based on a threshold [38]. Language models have also been used to fine-tune offline RL tasks [39] and to scaffold general sequential decision-making by initializing policies [40].

While prior works in HRL focus on skill discovery and option learning, and recent methods leverage LLMs for task decomposition and high-level planning, these approaches remain limited in their ability to automatically construct reusable, transferable skills for long-horizon tasks. Our work bridges these two directions by integrating LLM-guided subgoal discovery into the skill chaining framework, enabling the incremental construction of a hierarchy of generalizable options across diverse tasks.

III. PROBLEM FORMULATION

In this work, we address a general decision-making problem where the tasks for the agent are provided as natural language, denoted by $\{l_1, l_2, \dots, l_n\}$. Each l_i represents a task instruction. These task instructions are grounded through *goals*, defined as specific states the robot must reach to complete the task. The tasks are modeled as a semi-Markov Decision Process (semi-MDP) [14], defined by the tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, R, \tau, \gamma)$, where \mathcal{S} is the set of states, \mathcal{A} is the set of actions, $P(s'|s, a)$ is the probability of transitioning from state s to state s' after taking action a , $R(s, a)$ is the reward function, $\tau(s, a)$ is the time duration associated with action a , and γ is the discount factor. In this more general problem, the human language instructions $\{l_1, l_2, \dots, l_n\}$ specify high-level tasks, which may be complex or abstract. The robot should interpret these instructions and decompose them into a sequence of smaller subgoals that can be more easily achieved. Each subgoal corresponds to a set of states that satisfy the requirements of the subgoal within the semi-MDP, and the robot must develop a policy π that allows it

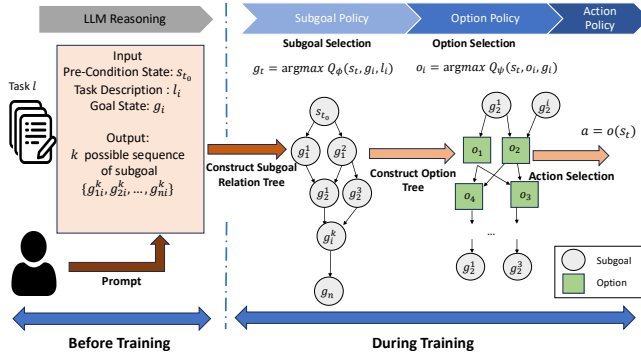


Fig. 2. Overview of the **LDSC** framework. The process consists of two phases: Before Training and During Training. In the Before Training phase, an LLM generates subgoal sequences based on task descriptions and initial conditions. During Training, a hierarchical structure operates with three components: (1) the subgoal policy, which constructs a subgoal relation tree and selects subgoals; (2) the option policy, which builds an option tree and determines the best option; and (3) the action policy, which selects the final action based on the chosen option.

to accomplish these tasks in a way that ultimately fulfills the overall objectives specified by $\{l_1, l_2, \dots, l_n\}$.

For a particular task l , the robot seeks to optimize its policy π to maximize the expected cumulative reward by following the optimal sequence of actions that lead to the completion of the main goal. Formally, the objective is to find the policy π^* that maximizes the expected discounted reward:

$$\pi^* = \arg \max_{\pi} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right] \quad (1)$$

where $\pi(a|s)$ is the policy that maps states to actions and $R(s_t, a_t)$ is the reward at time t . By breaking down complex goals into subgoals and solving them incrementally, the robot can handle more general tasks provided by natural language instructions.

IV. OUR APPROACH: LDSC

Our approach, LDSC, leverages reasoning to enhance learning efficiency by breaking tasks into manageable subgoals. This hierarchical structure enables the robot to learn a smaller, reusable set of options that can be applied across multiple subgoals. LDSC operates on three levels: the subgoal policy, which oversees high-level task planning and subgoal selection; the option policy, which operates at the intermediate level by selecting and executing the appropriate option based on the chosen subgoal; and the action policy, which handles detailed actions required to complete each subgoal.

A. LLM Reasoning

The algorithm begins with a human-provided instruction set $\{l_1, l_2, \dots, l_n\}$, which consists of a high-level goal, a description of the environment, constraints, and the initial state. Achieving such goals directly in sparse and complex environments is often difficult due to the large state space and the long horizons required to reach the goal. To address this challenge, the algorithm leverages the reasoning capabilities

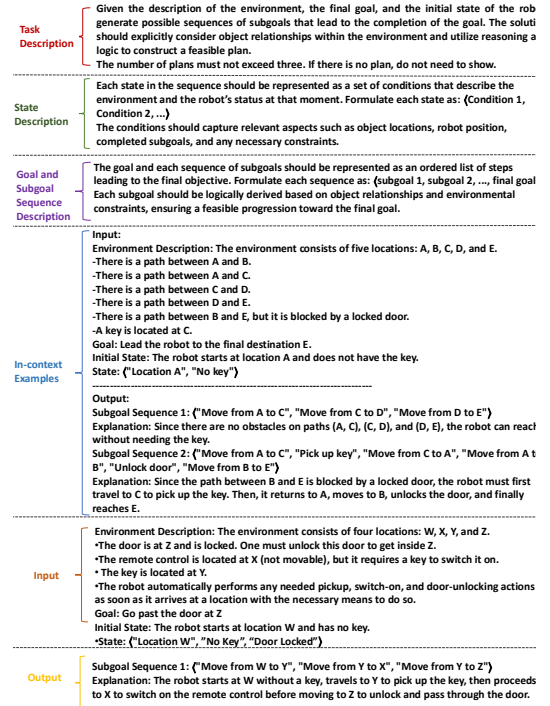


Fig. 3. Structured prompt design for generating subgoal sequences in a robotic planning task. The figure outlines different sections of the prompt, including task description, state representation, goal and subgoal sequencing, and in-context examples. The example shown serves as a qualitative example based on the Point Maze environment, demonstrating how a robot navigates the maze by reasoning over object relationships and environmental constraints to generate feasible and logically ordered action sequences.

of a LLM to decompose the given instruction set into a series of manageable subgoals $\{g_{11}, g_{12}, \dots, g_{mn}\}$. Formally, the LLM decouples the original goal into a sequence of subgoals:

$$G_n : \{g_{1i}, g_{2i}, \dots, g_{ji}\} = \text{LLM}(l_i, s_{t_0}) \quad (2)$$

where each g_{ji} represents a potential intermediate step j toward achieving the task i . This transformation allows the robot to focus on achieving smaller, more tractable goals in sequence, which are easier to accomplish than the original complex goal.

1) *Subgoal Relation Tree*: While LLMs possess strong reasoning capabilities, the subgoal sequences they suggest may not always be accurate or optimized. To address this, we draw inspiration from the Tree of Thought (ToT) framework [35] and propose constructing a subgoal relation tree to systematically explore and refine k possible subgoal sequences $\{G_n^1, G_n^2, \dots, G_n^k\}$ for task n . The tree is initialized with the current state s_{t_0} as the root. From this root, the LLM generates multiple subgoal sequences, where each sequence represents a potential path toward achieving the final task.

The generated subgoal sequences are directly connected to the root state s_{t_0} , forming the first layer of the subgoal relation tree. Specifically, each initial subgoal is defined as

$$g_{1i} = \text{create_child_subgoal}(s_{t_0}) \quad (3)$$

where g_{1i} represents the 1st subgoal for task i in the first layer. Subgoals that appear in multiple sequences are treated as shared nodes within the tree, enabling efficient representation of overlapping paths.

To expand the tree, for each subgoal g_{ji} in the current layer j , we generate a new set of subgoals for the next layer $j+1$ using the recursive relation

$$g_{(j+1)i} = \text{create_child_subgoal}(g_{ji}), \quad \text{if } (g_{ji}, g_{(j+1)i}) \notin E \quad (4)$$

where $g_{(j+1)i}$ represents a newly generated child subgoal of g_{ji} , and E denotes the set of existing edges in the tree. The expansion process continues recursively until one of the stopping criteria is met, such as a subgoal reaching the final task objective or the number of child nodes per subgoal exceeding a predefined threshold. This hierarchical structure enables systematic exploration of diverse task decompositions while ensuring that shared subgoals are efficiently reused across sequences. By leveraging this structured representation, planning efficiency is significantly improved for complex tasks.

B. Subgoal Policy

After constructing the relation tree, the algorithm proceeds to the training phase without further queries to the LLM. Next, the robot dynamically selects and evaluates goals based on the set of next goals in the relation tree. At each state s_t , the robot selects the goal g_t that maximizes its goal-value function. Since the robot must choose from a discrete set of goals, we employ Deep Q-learning (DQN) [3] to learn the goal-value function:

$$g_t = \arg \max_{g_i \in G} Q_\phi(s_t, g_i, l). \quad (5)$$

The goal-value function is updated using the Semi-Markov Decision Process (SMDP) Q-learning update [14]. Given an SMDP transition $(s_t, g_t, r_{t:t+\tau}, s_{t+\tau})$, the robot learns to improve its policy by optimizing the Q-value associated with selecting a goal g_t in state s_t . The target for updating the DQN, with Q-values parameterized by ϕ , is defined as:

$$y_t = \sum_{t'=t}^{t+\tau-1} \gamma^{t'-t} r_{t'} + \gamma^\tau Q_{\phi'}(s_{t+\tau}, \arg \max_{g' \in G} Q_\phi(s_{t+\tau}, g', l), l). \quad (6)$$

This formulation ensures that the robot can systematically select goals that maximize long-term rewards while accounting for the dynamics of the environment. By learning Q-values for state-goal pairs, the robot can prioritize goals that are most likely to achieve desired outcomes.

C. Option Policy

To lead the robot to the chosen subgoal, we choose o_t according to $\pi_O(s_t)$ using Equations 7 and 8. The option o_t is a temporally extended action that guides the robot through a sequence of actions to complete the subgoal. We can learn its option-value function using Deep Q-learning (DQN):

$$O'(s_t) = \{o_i \mid I_{o_i}(s_t) = 1 \cap \beta_{o_i}(s_t) = 0, \forall o_i \in O\} \quad (7)$$

$$o_t = \arg \max_{o_i \in O'(s_t)} Q_\psi(s_t, o_i, g) \quad (8)$$

Once the subgoals are defined, the robot employs a DSC approach to achieve each subgoal. At each step, the robot picks an option o_t based on its current policy $\pi_O(s_t, g_t)$. The environment is updated by executing the option, resulting in the reward $r_{t:t+\tau}$ and a new state $s_{t+\tau}$. After each option execution, the robot updates its policy using the collected data:

$$\pi_O = \text{update}(s_t, g_t, o_t, r_{t:t+\tau}, s_{t+\tau}) \quad (9)$$

The Q-value target for learning the weights ψ of the DQN is given by:

$$y_t = \sum_{t'=t}^{\tau} \gamma^{t'-t} r_{t'} + \gamma^{\tau-t} Q_{\psi'}(s_{t+\tau}, \arg \max_{o' \in O'(s_{t+\tau})} Q_\psi(s_{t+\tau}, o', g), g) \quad (10)$$

1) Option Tree Construction: The robot's behavior is structured through both the option tree and the subgoal relation tree. The subgoal relation tree is initialized with a starting subgoal at its root, and through this structure, the robot explores multiple possible paths to achieve the task. As the robot progresses through the subgoal tree, options are generated to handle transitions between subgoals.

Each subgoal in the relation tree corresponds to an option in the option tree. In particular, when the robot is at a given subgoal node g_i , options $\{o_1, \dots, o_k\}$ are created to transition between this subgoal and the next subgoal, represented as g_j . The structure of the option tree is recursively built as follows:

$$o^* = \text{create_child_option}(o) \quad (11)$$

where o^* represents a new option generated from the parent option o , which is used to connect two subgoals. The termination condition of each option is tied to the completion of its respective subgoal g_j , ensuring that the robot moves towards the next goal. As the robot progresses, the option tree grows, connecting various subgoals and forming a more intricate plan for task completion.

This hybrid tree structure allows the robot to connect subgoals using options, enabling it to break down complex tasks into manageable steps. The robot can effectively navigate through different levels of abstraction by leveraging the subgoal relation tree and option tree together.

D. Action Policy

The action policy is learned under the guidance of a high-level policy over options π_O . At the beginning of training, π_O contains only a single global option o_{z_i} for each subgoal i . The purpose of the global option is to explore the state space and reach the designated subgoal. This option has an initiation condition that holds across the entire state space $I_{o_{z_i}}(s) = 1$ for all s and terminates only upon reaching

Algorithm 1 LDSC

Given Human provides instruction set $L = \{l_1, l_2, \dots, l_n\}$, hyperparameter T_0 , the time budget of execution
Robot's Option Repertoire: $\mathcal{O} = \{\}$
Policy over tasks: $\pi_{\mathcal{O}'} : s_t \times l_i \rightarrow g_t$
Untrained Options: $o_U = \{\}$
Policy over options: $\pi_{\mathcal{O}} : s_t \times g_i \rightarrow o_t$
Buffer: $B = \{\}$

for all l_i in L **do**
 Use LLM to decouple instruction set into subgoals Set $\{G_i^1, G_i^2, \dots, G_i^k\}$ using Equation 2
 Construct the Subgoal Relation Tree using Equation 3
 for all new subgoal g_i **do**
 Global Option: $o_{z_i} = (I_{o_{z_i}}, \pi_{o_{z_i}}, \beta_{o_{z_i}} = 1_{g_i}, T = 1)$
 Goal Option: $o_{g_i} = (I_{o_{g_i}}, \pi_{o_{g_i}}, \beta_{o_{g_i}} = 1_{g_i}, T = T_0)$
 Add o_{g_i} into o_U
 end for
end for
for task l_i **do**
 $s_t = s_0$
 while t is not terminated **do**
 Choose a g using equation 5
 Choose a o using equation 8
 $r_{t:\tau}, s_{t+\tau} = \text{execute_option}(o_t)$
 $\pi_{\mathcal{O}} = \text{update}(s_t, g, o_t, r_{t:t+\tau}, s_{t+\tau})$ using Equation 10
 Buffer $B = B \cup r_{t:\tau}, s_{t+\tau}$
 if s_{t+1} meet g **then**
 load \hat{r}, \hat{s} in B
 $\pi_{\mathcal{O}'} = \text{update}(s_t, g, \hat{r}, \hat{s})$ using Equation 6
 Buffer $B = \{\}$
 end if
 if $\beta_{o_k}(s_{t+\tau})$ **and** $(s_0 \notin I_{o_i} \forall o_i \in \mathcal{O})$ **then**, $o_k \in o_U$
 $o_k.\text{learn_initiation_classifier}()$
 if $o_k.\text{initiation_classifier.is_trained}()$ **then**
 $\pi_{\mathcal{O}}.\text{add}(o_k)$
 $\mathcal{O}'.\text{append}(o_k)$
 Construct the Option Tree using Equation 11
 end if
 end if
 end while
end for

the subgoal $\beta_{o_{z_i}} = 1_{g_i}$. During execution, the global option selects primitive actions through its internal policy $\pi_{o_{z_i}}$ until termination.

After the global option o_{z_i} reaches the subgoal g_i a predefined number of times, a new option o_{g_i} is learned from the trajectories collected during these successful attempts. The goal option o_{g_i} shares the same termination condition as the global option, terminating upon reaching the subgoal state. Following this, a new option is established to reach the initiation set of the goal option o_{g_i} . By repeating this process, the set of available options in $\pi_{\mathcal{O}}$ is incrementally expanded as additional skills are discovered.

The learned options are executed over a predefined time horizon T . When the robot selects an option o in state s_t , the option executes its closed-loop control policy for τ time steps until either its termination condition is satisfied or the time limit T is reached. Once the option terminates or times out, control is returned to $\pi_{\mathcal{O}}$, which selects the next option to execute from the resulting state $s_{t+\tau}$.

TABLE I

HYPERPARAMETERS USED IN THE EXPERIMENTS.

Parameter	Value
Batch size (N)	64
Discount factor (γ)	0.99
Soft update coefficient (τ)	0.01
Hidden sizes (DDPG)	[400, 300]
Hidden layers (DQN)	[32, 32]
Critic learning rate	1×10^{-3}
Actor learning rate	1×10^{-4}

V. EXPERIMENTS

We evaluate our algorithm in Mujoco [41] framework for RL research. Mujoco is a framework for physical dynamic simulation, often employed to test HRL algorithms. We evaluate our algorithm on four distinct map settings. We apply ChatGPT-o1 for the reasoning. LLM reasoning Details of the network hyperparameters are provided in Table I.

A. Mujoco

Four Rooms with Lock and Key: This map is adapted from the Deep Skill Chaining paper [13], which extends the original Four Rooms environment by introducing a key and lock mechanism. A point robot [42] is placed in the Four Rooms setting, where its objective is to first pick up the key and then navigate to the lock, represented by a red sphere in the top-left room. The robot's state space includes its position, orientation, linear velocity, rotational velocity, and a binary "has-key" indicator. The robot must follow the LLM output sequence: `<"retrieve key", "reach the lock">`. If the robot reaches the lock while holding the key, the episode terminates with a sparse reward of 1. Otherwise, the robot incurs a step penalty.

Point Maze: We extend the Point Maze environment by introducing two subgoal checkpoints. The robot must navigate through these checkpoints sequentially to receive rewards. The LLM output sequence follows the order: `<"retrieve key (green ball)", "Switch on the Remote Control (blue ball)", "go to the door (red ball)">`. The state space and reward structure are designed to emphasize long-term planning and subgoal discovery. The episode terminates once all goals have been traversed. The prompt structure of this task is illustrated in Figure 3.

Point E-Maze: We extend the benchmark U-shaped Point-Maze task [13] by introducing two possible starting positions for the robot: one at the top and one at the bottom rungs of the E-shaped maze. Furthermore, we modify the task by incorporating two subgoals in the form of keys, located at the top-right and bottom-right sections of the maze. To reach the final goal, the robot must first collect both keys, requiring it to exhibit strategic planning and skill chaining. The robot can follow either of the two possible LLM outputs: `<"retrieve key1 (blue ball)", "retrieve key2 (green ball)", "reach the goal (red ball)">` or `<"retrieve key2 (green ball)",`

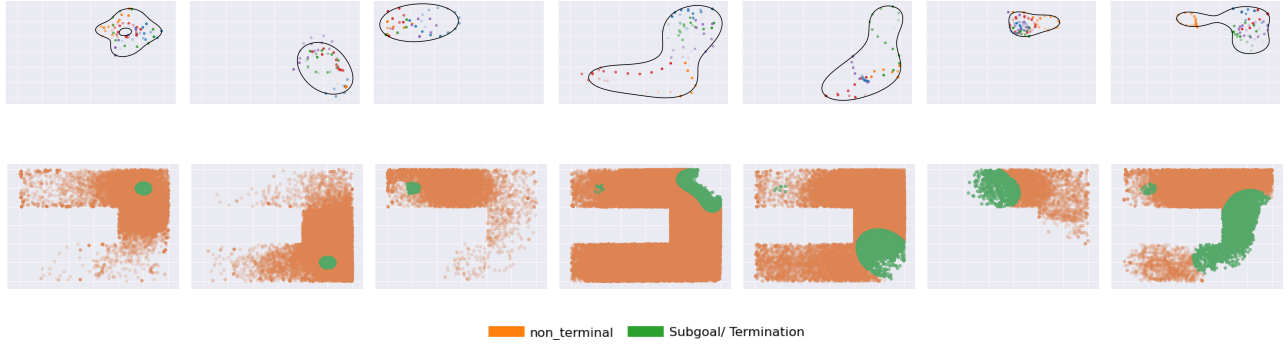


Fig. 4. **Qualitative performance** of the robot in the Point Maze environment. The upper row shows the initial set for each option, illustrating the state space coverage where the option can be executed. The lower row displays the corresponding policy plots, where orange regions indicate areas where the policy continues execution, while green regions signify termination states. The robot follows a structured sequence: first reaching subgoal 1 (top-right), then subgoal 2 (bottom-right), and finally the goal (top-left).

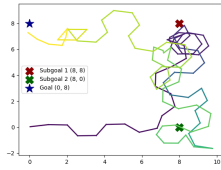


Fig. 5. Trajectory of the quantitative example, where the path is colored indicating the change of options. Each segment of the trajectory represents the agent's movement in 2D space while executing a specific option.

"retrieve key1 (blue ball)", "reach the goal (red ball)">. This extension challenges the robot to efficiently coordinate subgoal completion before progressing to the final objective.

Tunnel: The Tunnel environment presents a unique challenge in robot navigation within constrained spaces, evolving the traditional maze by incorporating long, narrow tunnels. These tunnels split into two distinct branches, with a checkpoint located before the branches. At the entrance, the robot must first collect the green key (key1) to unlock the path to the goal. The environment tests the robot's ability to plan and sequence actions strategically. The robot can follow one of the following possible LLM output sequences to complete the task: <"go to the checkpoint", "retrieve key1 (green)", "proceed to the goal">, or <"go to the checkpoint", "retrieve key1 (green)", "collect key2 (red)", "proceed to the goal">, or <"go to the checkpoint", "retrieve key2 (red)", "collect key1 (green)", "proceed to the goal">. These variations provide the robot with different strategic options, requiring high-level planning.

B. Baselines

We have considered the following baseline algorithms. **Deep Skill Chaining (DSC)** [13]: A HRL method that decomposes complex tasks into a sequence of sub-skills, enabling more efficient learning and planning. **Deep Deterministic Policy Gradient (DDPG)** [43]: A model-free, off-policy actor-critic algorithm designed for continuous action spaces,

leveraging experience replay and target networks for stable training. **Option-Critic** [20]: A framework that combines the option-based hierarchy with a policy gradient method, allowing the robot to learn both policies over options and policies within options.

C. Evaluation

1) *Qualitative Example:* To assess the qualitative performance of the proposed LDSC framework, we present visual examples that illustrate its behavior in a Maze environment. As shown in Figure 4, LDSC effectively learns high-level abstract options that enable structured and hierarchical navigation toward subgoals and the final goal. The first three options demonstrate localized initiation sets and policies that guide the robot to key subgoal locations, including the top-right, top-left, and bottom-right corners of the maze. Beyond these localized strategies, LDSC also generates bridging options that facilitate transitions between subgoals and the final goal. Notably, as shown in Figure 5, the trajectory of the quantitative example demonstrates how the agent transitions between different options, with the path colored according to the active option, demonstrating its capability to form temporally extended and goal-directed policies.

This hierarchical decomposition of the navigation task into smaller, modular subproblems significantly improves the efficiency of planning and execution. By leveraging learned options, the agent effectively reduces unnecessary exploration, ensuring more directed and optimized movement between subgoals and the final destination. The ability to abstract and generalize decision-making at different levels of granularity further highlights the strength of LDSC in solving complex, long-horizon navigation tasks.

2) *Quantitative Comparison:* We evaluate the robot's performance by assessing its task completion rates and the time taken to complete each task based on human-provided goals. The corresponding visual representations are shown in Figures 6 (Upper), while the performance results are depicted in Figures 6 (Bottom).

Our approach, LDSC, demonstrates superior performance compared to all baseline algorithms without requiring additional training time, as the LLM processing is completed

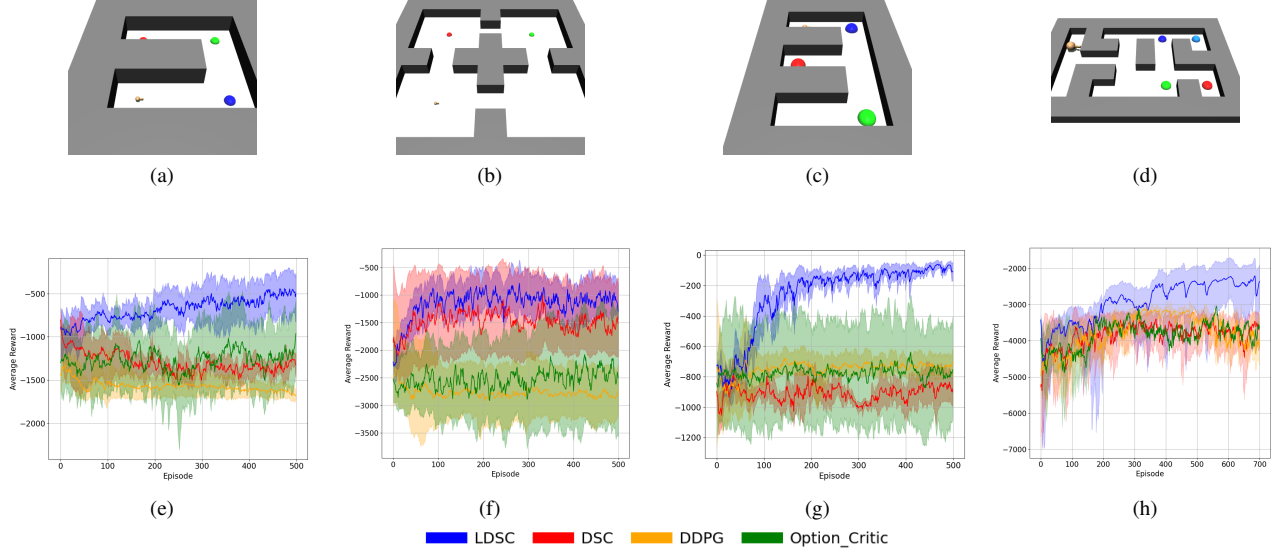


Fig. 6. (Upper) Environment setup and (Bottom) average reward for each task: (a) and (e) are Maze, (b) and (f) are Four Rooms, (c) and (g) are E-Maze, and (d) and (h) are Tunnel.

Map	DSC		DDPG		Option-Critic		LDSC (Ours)	
	Success Rate \uparrow	Time (s) \downarrow	Success Rate \uparrow	Time (s) \downarrow	Success Rate \uparrow	Time (s) \downarrow	Success Rate \uparrow	Time (s) \downarrow
Maze	0% \pm 0%	1063 \pm 274	0% \pm 0%	1187 \pm 116	0% \pm 0%	1030 \pm 147	100% \pm 0%	485 \pm 174
FourRoom	86% \pm 8%	1035 \pm 479	0% \pm 0%	1331 \pm 355	0% \pm 0%	1181 \pm 500	95% \pm 2%	678 \pm 208
E-Maze	0% \pm 0%	1345 \pm 145	0% \pm 0%	1344 \pm 83	0% \pm 0%	960 \pm 306	100% \pm 0%	86.5 \pm 12.5
Tunnel	0% \pm 0%	1368 \pm 370	0% \pm 0%	1527 \pm 527	0% \pm 0%	1349 \pm 349	81.8% \pm 3.6%	906 \pm 306

TABLE II

COMPARISON OF DIFFERENT METHODS IN TERMS OF SUCCESS RATE AND COMPLETION TIME ACROSS VARIOUS MAPS.

beforehand. Notably, while DSC performs better than the baseline methods, LDSC achieves a significant margin of improvement across all tasks. On average, LDSC outperforms baseline methods by **55.9%**, reduces task completion time by **53.1%**, and improves success rates by **72.7%**, ensuring more efficient and reliable execution.

The key advantage of LDSC lies in its LLM-driven hierarchical decomposition, where the problem is broken down into smaller, localized subproblems using reasoning capabilities from a language model. By leveraging LLM-based reasoning, LDSC can infer logical subgoal structures, effectively segmenting the task into meaningful intermediate steps. This hierarchical structure significantly reduces the search space, thereby enhancing exploration efficiency. Since the agent can now focus on reaching subgoals rather than exploring the entire state space at once, it requires fewer interactions to discover optimal behaviors. As a result, the learning speed is significantly accelerated, and the agent exhibits a more robust and stable performance across different maze configurations. This is particularly evident in the Point E-Maze setting, where the introduction of multiple subgoals increases task complexity, yet LDSC still outperforms other methods by a substantial margin. We also evaluate the robot's performance in a complex environment, Tunnel, which requires long-term planning. LDSC significantly outperforms

baseline methods, demonstrating superior efficiency and adaptability in challenging scenarios.

These results highlight the effectiveness of LDSC in leveraging high-level abstraction and LLM-driven reasoning to improve long-horizon planning, ultimately leading to higher task success rates and more efficient navigation in complex environments.

VI. CONCLUSION

In this work, we introduced LDSC, a semantic HRL method that leverages LLMs for subgoal reasoning and decomposition. Our approach effectively addresses key challenges in exploration efficiency, policy generalization, and option reusability in complex, multi-task reinforcement learning settings. By integrating LLM-driven reasoning, LDSC constructs structured subgoal hierarchies, significantly improving task decomposition and learning efficiency. This is achieved without increasing overall training time, as the LLM operates only during the pre-processing stage. Through extensive experiments in diverse maze environments, we demonstrated that LDSC outperforms the existing RL approaches DSC, DDPG, and Option-Critic methods. Our results show that LDSC achieves a substantial improvement in task success rates, task completion time, and overall learning efficiency, highlighting its robustness and adaptability across different

environments.

The significance of LDSC extends beyond the current experiments. Future research can explore scaling LDSC to high-dimensional robotic tasks and further refining semantic reasoning mechanisms to enhance adaptability. Additionally, investigating multi-agent collaboration and lifelong learning extensions for LDSC could open new directions in RL.

Overall, LDSC provides a compelling framework for hierarchical decision-making, offering a scalable and interpretable approach to RL in dynamic and uncertain environments.

REFERENCES

- [1] J. Liu, P. Hang, X. Qi, J. Wang, and J. Sun, "Mtd-gpt: A multi-task decision-making gpt model for autonomous driving at unsignalized intersections," in *2023 IEEE 26th International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2023, pp. 5154–5161.
- [2] J. Fu, Y. Long, K. Chen, W. Wei, and Q. Dou, "Multi-objective cross-task learning via goal-conditioned gpt-based decision transformers for surgical robot task automation," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024, pp. 13 362–13 368.
- [3] V. Mnih, "Playing atari with deep reinforcement learning," *arXiv preprint arXiv:1312.5602*, 2013.
- [4] Y. Burda, H. Edwards, A. Storkey, and O. Klimov, "Exploration by random network distillation," *arXiv preprint arXiv:1810.12894*, 2018.
- [5] J. Schulman, "Trust region policy optimization," *arXiv preprint arXiv:1502.05477*, 2015.
- [6] P. Henderson, R. Islam, P. Bachman, J. Pineau, D. Precup, and D. Meger, "Deep reinforcement learning that matters," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 32, no. 1, 2018.
- [7] Y. Tao, S. Genc, J. Chung, T. Sun, and S. Mallya, "Repaint: Knowledge transfer in deep reinforcement learning," in *International conference on machine learning*. PMLR, 2021, pp. 10 141–10 152.
- [8] V. Dhiman, S. Banerjee, B. Griffin, J. M. Siskind, and J. J. Corso, "A critical investigation of deep reinforcement learning for navigation," *arXiv preprint arXiv:1802.02274*, 2018.
- [9] A. A. Rusu, S. G. Colmenarejo, C. Gulcehre, G. Desjardins, J. Kirkpatrick, R. Pascanu, V. Mnih, K. Kavukcuoglu, and R. Hadsell, "Policy distillation," *arXiv preprint arXiv:1511.06295*, 2015.
- [10] A. G. Barto and S. Mahadevan, "Recent advances in hierarchical reinforcement learning," *Discrete event dynamic systems*, vol. 13, pp. 341–379, 2003.
- [11] R. S. Sutton, D. Precup, and S. Singh, "Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning," *Artificial intelligence*, vol. 112, no. 1-2, pp. 181–211, 1999.
- [12] S. Thrun and A. Schwartz, "Finding structure in reinforcement learning," *Advances in neural information processing systems*, vol. 7, 1994.
- [13] A. Bagaria and G. Konidaris, "Option discovery using deep skill chaining," in *International Conference on Learning Representations*, 2019.
- [14] S. Bradtke and M. Duff, "Reinforcement learning methods for continuous-time markov decision problems," in *Advances in Neural Information Processing Systems*, G. Tesauro, D. Touretzky, and T. Leen, Eds., vol. 7. MIT Press, 1994. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/1994/file/07871915a8107172b3b5dc15a6574ad3-Paper.pdf
- [15] G. Konidaris and A. Barto, "Skill discovery in continuous reinforcement learning domains using skill chaining," *Advances in neural information processing systems*, vol. 22, 2009.
- [16] M. C. Machado, C. Rosenbaum, X. Guo, M. Liu, G. Tesauro, and M. Campbell, "Eigenoption discovery through the deep successor representation," *arXiv preprint arXiv:1710.11089*, 2017.
- [17] B. Liu, Y. Jiang, X. Zhang, Q. Liu, S. Zhang, J. Biswas, and P. Stone, "Llm+ p: Empowering large language models with optimal planning proficiency," *arXiv preprint arXiv:2304.11477*, 2023.
- [18] M. Kwon, S. M. Xie, K. Bullard, and D. Sadigh, "Reward design with language models," *arXiv preprint arXiv:2303.00001*, 2023.
- [19] A. Sharma, S. Gu, S. Levine, V. Kumar, and K. Hausman, "Dynamics-aware unsupervised discovery of skills," *arXiv preprint arXiv:1907.01657*, 2019.
- [20] P.-L. Bacon, J. Harb, and D. Precup, "The option-critic architecture," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 31, no. 1, 2017.
- [21] R. Chunduru and D. Precup, "Attention option-critic," *arXiv preprint arXiv:2201.02628*, 2022.
- [22] R. Sutton, D. Precup, and S. Singh, "Intra-option learning about temporally abstract actions," 01 1998, pp. 556–564.
- [23] R. Ramesh, M. Tomar, and B. Ravindran, "Successor options: An option discovery framework for reinforcement learning," 2019. [Online]. Available: <https://arxiv.org/abs/1905.05731>
- [24] R. Fox, S. Krishnan, I. Stoica, and K. Goldberg, "Multi-level discovery of deep options," *arXiv preprint arXiv:1703.08294*, 2017.
- [25] B. Eysenbach, A. Gupta, J. Ibarz, and S. Levine, "Diversity is all you need: Learning skills without a reward function," *arXiv preprint arXiv:1802.06070*, 2018.
- [26] R. Ghugare, M. Geist, G. Berseth, and B. Eysenbach, "Closing the gap between td learning and supervised learning—a generalisation point of view," *arXiv preprint arXiv:2401.11237*, 2024.
- [27] W. Huang, P. Abbeel, D. Pathak, and I. Mordatch, "Language models as zero-shot planners: Extracting actionable knowledge for embodied agents," in *International conference on machine learning*. PMLR, 2022, pp. 9118–9147.
- [28] D. Zhou, N. Schärli, L. Hou, J. Wei, N. Scales, X. Wang, D. Schuurmans, C. Cui, O. Bousquet, Q. Le, et al., "Least-to-most prompting enables complex reasoning in large language models," *arXiv preprint arXiv:2205.10625*, 2022.
- [29] C. Pan, X. Yang, H. Wang, W. Wei, and T. Li, "Hierarchical continual reinforcement learning via large language model," *arXiv preprint arXiv:2401.15098*, 2024.
- [30] C. H. Song, J. Wu, C. Washington, B. M. Sadler, W.-L. Chao, and Y. Su, "Llm-planner: Few-shot grounded planning for embodied agents with large language models," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 2998–3009.
- [31] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al., "Chain-of-thought prompting elicits reasoning in large language models," *Advances in neural information processing systems*, vol. 35, pp. 24 824–24 837, 2022.
- [32] L. Wang, W. Xu, Y. Lan, Z. Hu, Y. Lan, R. K.-W. Lee, and E.-P. Lim, "Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models," *arXiv preprint arXiv:2305.04091*, 2023.
- [33] A. Zhou, K. Yan, M. Shlapentokh-Rothman, H. Wang, and Y.-X. Wang, "Language agent tree search unifies reasoning acting and planning in language models," *arXiv preprint arXiv:2310.04406*, 2023.
- [34] M. Hu, Y. Mu, X. Yu, M. Ding, S. Wu, W. Shao, Q. Chen, B. Wang, Y. Qiao, and P. Luo, "Tree-planner: Efficient close-loop task planning with large language models," *arXiv preprint arXiv:2310.08582*, 2023.
- [35] S. Yao, D. Yu, J. Zhao, I. Shafraan, T. Griffiths, Y. Cao, and K. Narasimhan, "Tree of thoughts: Deliberate problem solving with large language models," *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [36] W. Yu, N. Gileadi, C. Fu, S. Kirmani, K.-H. Lee, M. Gonzalez Arenas, H.-T. Lewis Chiang, T. Erez, L. Hasenclever, J. Humplik, B. Ichter, T. Xiao, P. Xu, A. Zeng, T. Zhang, N. Heess, D. Sadigh, J. Tan, Y. Tassa, and F. Xia, "Language to rewards for robotic skill synthesis," *arXiv preprint arXiv:2306.08647*, 2023.
- [37] H. Li, X. Yang, Z. Wang, X. Zhu, J. Zhou, Y. Qiao, X. Wang, H. Li, L. Lu, and J. Dai, "Auto mc-reward: Automated dense reward design with large language models for minecraft," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2024, pp. 16 426–16 435.
- [38] J. Rocamonde, V. Montesinos, E. Nava, E. Perez, and D. Lindner, "Vision-language models are zero-shot reward models for reinforcement learning," *arXiv preprint arXiv:2310.12921*, 2023.
- [39] M. Reid, Y. Yamada, and S. S. Gu, "Can wikipedia help offline reinforcement learning?" *arXiv preprint arXiv:2201.12122*, 2022.
- [40] S. Li, X. Puig, C. Paxton, Y. Du, C. Wang, L. Fan, T. Chen, D.-A. Huang, E. Akyürek, A. Anandkumar, et al., "Pre-trained language models for interactive decision-making," *Advances in Neural Information Processing Systems*, vol. 35, pp. 31 199–31 212, 2022.
- [41] E. Todorov, T. Erez, and Y. Tassa, "Mujoco: A physics engine for model-based control," in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2012, pp. 5026–5033.
- [42] Y. Duan, X. Chen, R. Houthoofd, J. Schulman, and P. Abbeel, "Benchmarking deep reinforcement learning for continuous control," in *International conference on machine learning*. PMLR, 2016, pp. 1329–1338.

- [43] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, 2015.