

Deconstructing Spatial Complexity: Hierarchical Decomposition for LLM Spatial Reasoning

Anonymous ACL submission

Abstract

LLMs have shown remarkable proficiency in general language understanding and reasoning. However, they consistently underperform in spatial reasoning that severely limits their application, particularly in embodied intelligence. Inspired by the success of hierarchical reinforcement learning, this paper introduces a novel method for hierarchical task decomposition in LLM spatial reasoning. Our approach guides LLMs to decompose complex tasks into manageable sub-tasks by identifying key intermediate states and generating simplified sub-environments. However, we identify that LLMs often fail to derive optimal intermediate states due to their insufficient spatial prior, leading to sub-optimal task decomposition. To address this limitation and enhance its planning capability, we propose the MCTS-Guided Group Relative Policy Optimization (M-GRPO), where we reformulate the UCT formula by incorporating the LLM’s prior predictive probabilities alongside its epistemic uncertainty. Furthermore, we implement a more fine-grained advantage function, enabling the model to learn optimal path planning. Experimental results demonstrate that our method substantially improves LLM performance on spatial tasks, including navigation, planning, and strategic games, achieving state-of-the-art results. This work paves the way for LLMs in real-world applications.

1 Introduction

Large Language Models (LLMs) have revolutionized the landscape of artificial intelligence, achieving remarkable breakthroughs across various domains, including natural language processing and scientific reasoning (Zhao et al., 2023). However, as LLMs transition into the era of embodied AI, a critical and persistent bottleneck has emerged: their inherent limitations in spatial reasoning. While LLMs excel at manipulating abstract concepts and language, they often struggle with understanding

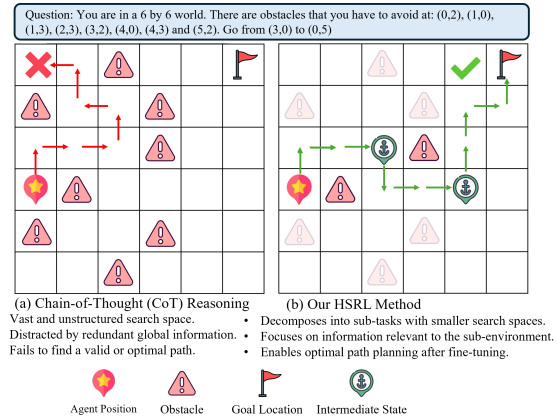


Figure 1: Comparison between CoT Reasoning and our HSRL method. (a) Standard CoT reasoning fails on complex spatial planning by inefficiently exploring a vast search space amidst distracting information. (b) In contrast, our HSRL framework succeeds by decomposing the task via key intermediate states and constructing focused sub-environments, which enables efficient and optimal planning.

complex spatial relationships, performing efficient path planning, and engaging in sequential action reasoning (Ma et al., 2025; Chen et al., 2024). This severely limits their development and practical deployment in embodied systems.

Existing research has explored several avenues to address this challenge, yet each faces significant limitations. Prompt engineering methods like CoT (Wei et al., 2022b), ToT (Yao et al., 2023a) and ProgPrompt (Singh et al., 2023) aim to elicit reasoning through specialized prompts, but their effectiveness is capped by the model’s often flawed intrinsic spatial capabilities. Fine-tuning approaches (Dao and Vu, 2025; Deng et al., 2025; Aghzal et al., 2024b) show promise but typically demand vast, expensive task-specific datasets and suffer from poor generalization to novel environments. Task decomposition strategies like HyperTree (Gui et al., 2025) and Plan-and-Act (Erdogan et al., 2025b) are primarily designed for tasks with clear, language-

063 based "logical breaks", rendering them ill-suited
064 for spatial reasoning problems like pathfinding that
065 lack such linguistic segmentation. Finally, offload-
066 ing planning to external, non-differentiable tools
067 breaks the end-to-end optimization paradigm, as
068 these tools cannot be jointly trained with the LLM's
069 representation layer and may not be universally de-
070 ployable at test time.

071 To overcome these limitations, we introduce Hi-
072 erarchical Spatial Reasoning with LLM (HSRL),
073 a novel hierarchical spatial reasoning paradigm
074 inspired by hierarchical reinforcement learning.
075 The core innovation of HSRL lies in its state- and
076 environment-based hierarchical mechanism, which
077 is fundamentally different from prior language-
078 based decomposition methods. The framework
079 employs a two-level hierarchy. A high-level LLM
080 planner decomposes complex tasks by generating
081 a structured sequence of intermediate states, which
082 serve as sub-goals to guide the overall planning
083 process. Then, for each state-to-state transition, a
084 low-level LLM agent performs a dual role, where it
085 first acts as an environment processor to construct a
086 localized sub-environment, and subsequently oper-
087 ates as an actor to execute precise actions required
088 to reach the target sub-goal. While this hierarchical
089 structure provides a powerful framework for de-
090 composition, the high-level planner's effectiveness
091 is constrained by the LLM's insufficient spatial
092 priors, leading to the generation of sub-optimal in-
093 termediate states. To address this, we propose an in-
094 novative online fine-tuning framework, M-GRPO,
095 designed to enhance the high-level planner. Our
096 approach improves planning optimality by tackling
097 two fundamental challenges: effective exploration
098 of the solution space and precise credit assignment
099 for training. To achieve robust exploration, we
100 draw inspiration from Monte Carlo Tree Search
101 (MCTS), where the high-level LLM generates mul-
102 tiple candidate sequences of intermediate states,
103 building a search tree to systematically explore
104 diverse planning strategies. Unlike standard LLM-
105 MCTS integration paradigms (Xie et al.; Guan
106 et al.) rely on the vanilla UCT formula, we pro-
107 pose incorporating the LLM's intrinsic confidence
108 and uncertainty as internal reward signals. This
109 integration facilitates the selection of trajectories
110 that exhibit both high environmental rewards and
111 high internal certainty, thereby mitigating the in-
112 stability caused by high-reward but low-confidence
113 paths—which often represent stochastic hallucina-
114 tions that hinder convergence.(Kuhn et al., 2023;

Azaria and Mitchell, 2023) Furthermore, we lever-
age perplexity-based uncertainty as an informative
metric to evaluate whether the model has suffi-
ciently explored the state space within its semantic
knowledge, leading to a more principled and effi-
cient search process. For precise credit assignment,
we introduce a fine-grained advantage function, a
significant departure from traditional Group Rela-
tive Policy Optimization (GRPO) which evaluates
whole-trajectory values without detailed supervi-
sion. Our method calculates the advantage of each
intermediate state relative to its "sibling" states (i.e.,
those that share a common prefix state sequence).
This provides a focused and accurate training sig-
nal, enabling the LLM to learn which specific sub-
goals are the most effective. Our method requires
only a small amount of data and can be flexibly
applied to multi-level planning tasks. In summary,
this work makes the following key contributions:

- **A Novel State- and Environment-Based Hierarchical Reasoning Framework:** We introduce HSRL, a framework that presents a novel state- and environment-based decomposition paradigm for LLM spatial reasoning, departing from prevalent language-based methods. This paradigm is specifically designed to address continuous spatial problems where traditional language-based decomposition is ineffective.
- **A Novel Fine-Tuning Framework for Planning Optimality:** To address the sub-optimal planning inherent in pre-trained LLMs, we develop M-GRPO, a new fine-tuning algorithm. By integrating a modified Monte Carlo Tree Search exploration mechanism with a fine-grained, node-level advantage function, our method substantially improves planning optimality with high data efficiency.
- **Comprehensive Empirical Validation of Superiority:** Through extensive experiments on large-scale navigation, object planning, and strategy game benchmarks, we demonstrate that HSRL achieves state-of-the-art performance. The results validate its significant gains over existing methods and its strong generalization across diverse task modalities.

2 Method

In this work, we introduce the HSRL framework, as illustrated in Figure 2, to address the limitations

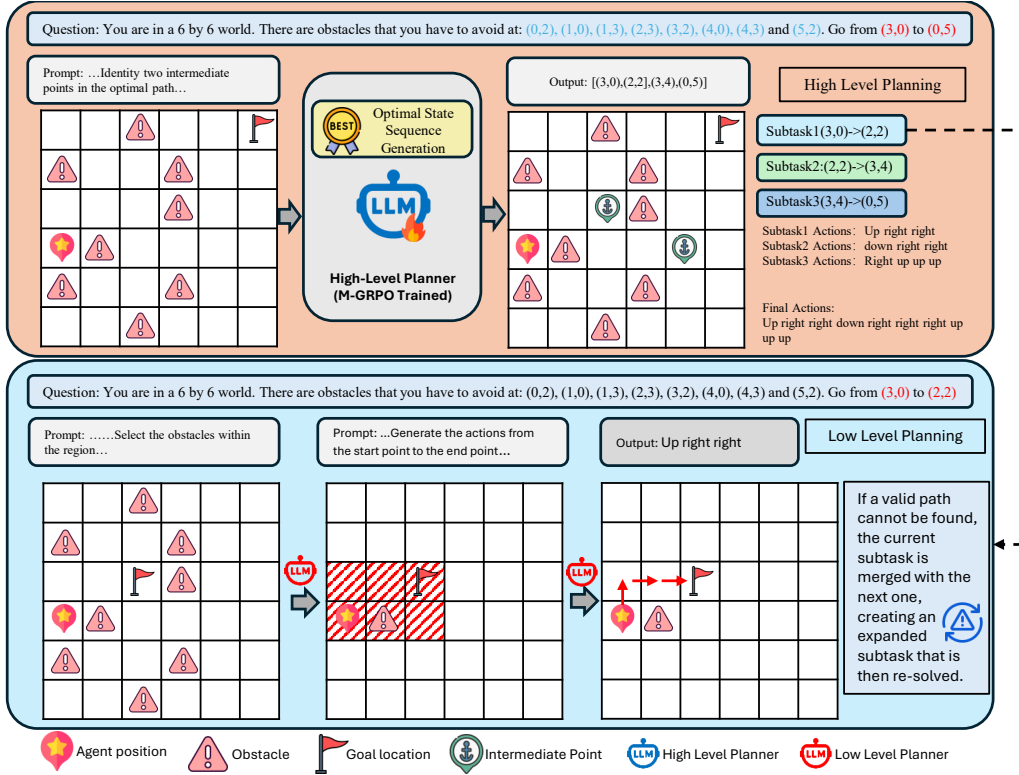


Figure 2: An overview of the HSRL framework. This framework employs a two-level hierarchical strategy. An M-GRPO trained high-level planner first identifies key intermediate states, decomposing the task into a series of sub-tasks. A low-level planner selects relevant information for the sub-task and then generates action sequences for each sub-task within a localized sub-environment. If a sub-task is unsolvable, it is merged with the next one (e.g., from the start of Subtask 1 to the end of Subtask 2) and replanned.

of existing LLM-based planning methods (see Appendix A). Our approach consists of two key components: a novel two-level hierarchical framework that decomposes complex tasks into a series of manageable sub-problems, and an innovative MCTS-guided finetuning method designed to enhance the optimality of the generated plans.

2.1 Hierarchical Planning with State and Environment Decomposition

Our framework leverages a two-level hierarchical decomposition strategy to break down complex planning tasks. This decomposition is applied at both the state level and the environmental level, effectively managing the complexity of the problem space.

State-Level Decomposition via LLM. Prior research in LLM-based path planning has shown promising results by manually decomposing tasks into sub-goals (Aghzal et al., 2024b). We extend this concept by enabling the LLM to autonomously generate these key intermediate states. Given a task’s initial and final states, our method prompts

the LLM to reason and generate a concise sequence of critical intermediate states. This process transforms a high-level goal into a series of state-to-state transitions, effectively simplifying the planning horizon for subsequent steps.

Environmental-Level Decomposition and Dynamic Expansion. After decomposing the task at the state level, much of the global environmental information becomes irrelevant noise for solving a specific sub-task, which can hinder the reasoning process. Following the generation of the state sequence, we define a sub-task for each consecutive pair of intermediate states. For each sub-task, we create a corresponding regional environment by identifying information that is relevant to the current sub-problem (e.g., obstacles or landmarks within a localized area). This hierarchical representation allows the model to focus on a smaller, more manageable sub-environment, thereby improving efficiency and reducing the search space. If the model is unable to find a valid path within the localized environment, the scope of the sub-task is expanded. The end state is extended to the next

intermediate state in the sequence, creating a larger sub-task that encompasses a broader area. This process is repeated until a solution is found or, in the worst case, the problem reverts to the original, full-scale task, ensuring robust and complete coverage of the problem space.

2.2 Optimizing Planning with M-GRPO

Due to an inherent lack of pre-training in spatial reasoning, LLMs often struggle to generate optimal sequences of intermediate states. However, the generation of correct intermediate states is a critical prerequisite for the success of subsequent environment decomposition and low-level action planning. To address this, we propose an online learning approach, as illustrated in Figure 3, that integrates the exploratory power of MCTS with the fine-tuning process of GRPO. This approach enables the LLM to learn and improve its planning policy during exploration.

MCTS-Guided Exploration for Optimal State

Generation The state generation process is framed as a search problem navigated by MCTS. Within this search tree, each node represents an intermediate state, while a full sequence of nodes forms a complete trajectory, known as a completion. Starting from the initial state, the tree is built iteratively. In each iteration, the MCTS policy traverses the tree to a leaf node. From this state, the LLM is prompted to generate subsequent potential states (expansion), a reward is evaluated (simulation), and the Q-values along the path are updated (backpropagation). To effectively integrate the LLM’s prior knowledge with empirical feedback during the search process, we propose a refined UCT heuristic that dynamically tunes the selection through a confidence weight $c(s')$ and an uncertainty factor $u(s')$. Specifically, based on the average log-probability $\ell(s')$ of the generated sequence

$$\ell(s') = \frac{1}{|T|} \sum_{t=1}^{|T|} \log P(x_t | x_{<t}, \text{context}) \quad (1)$$

where T denotes the set of tokens in the generated intermediate state, and $P(x_t | x_{<t}, \text{context})$ represents the likelihood of each token assigned by the LLM. We define the confidence weight as:

$$c(s') = \exp(\tau \cdot \ell(s')) \quad (2)$$

which calibrates the exploitation term to ensure the search prioritizes semantically coherent paths.

Algorithm 1 M-GRPO Training Algorithm

Require: Policy π_θ , initial state s_0 , iterations N_{\max} , group size M
Ensure: Optimized policy π_θ
1: $T \leftarrow \text{InitializeTree}(s_0)$
2: **for** it = 1 **to** N_{\max} **do**
3: // MCTS-Guided Exploration
4: $L \leftarrow \text{SELECTLEAF}(T)$ using Refined UCT (Eq. 4)
5: $\{\tau_m\}_{m=1}^M \leftarrow \text{EXPANDANDSIMULATE}(\pi_\theta, L)$
6: **for** each trajectory τ_m **do**
7: $R_m \leftarrow \text{SIMULATETOGOAL}(\tau_m)$
8: $\text{BACKPROPAGATE}(\tau_m, R_m)$ {Update Q-values}
9: **end for**
10: // Fine-Grained Advantage Estimation
11: $A_{\text{all}} \leftarrow []$
12: **for** each $\tau_m = \{s_{m,1}, \dots, s_{m,T}\}$ **do**
13: **for** $n = 1$ **to** T **do**
14: $\mathcal{S}_{\text{sib}} \leftarrow \text{GETSIBLINGS}(s_{m,n})$
15: $\bar{Q}_{\text{sib}} \leftarrow \text{MEAN}\{Q(s') \mid s' \in \mathcal{S}_{\text{sib}}\}$
16: $A_{m,n} \leftarrow Q(s_{m,n}) - \bar{Q}_{\text{sib}}$ {Node-level adv.}
17: **end for**
18: $A(\tau_m) \leftarrow \frac{1}{T} \sum_{n=1}^T A_{m,n}$ {Trajectory adv.}
19: $A_{\text{all}}.\text{APPEND}(A(\tau_m))$
20: **end for**
21: $\pi_\theta \leftarrow \text{UPDATEPOLICY}(\pi_\theta, A_{\text{all}})$ {GRPO Loss Update}
22: **end for**
23: **return** π_θ

Simultaneously, an uncertainty factor is defined as:

$$u(s') = 1 + \gamma \cdot \text{clamp}(-\ell(s'), 0, u_{\max}) \quad (3)$$

which utilizes log-perplexity as a gain for the exploration term to guide the algorithm toward the model’s epistemic blind spots. In cases where multiple candidate sequences share the same intermediate state s' , we average their respective $c(s')$ and $u(s')$ to ensure a representative estimation. The final selection policy is formulated as:

$$s_{\text{next}} = \underset{s' \in \text{Children}(s)}{\text{argmax}} \left\{ Q(s') \cdot c(s') + C \cdot u(s') \sqrt{\frac{\ln N(s)}{N(s')}} \right\} \quad (4)$$

where $Q(s')$ represents the estimated value, $N(s)$ and $N(s')$ denote visit counts, and C is the exploration constant. This design establishes a dynamic mechanism between model intuition and empirical value, significantly enhancing exploration efficiency.

Fine-Grained Advantage Function for Precise Policy Updates

In standard policy optimization frameworks like GRPO, the advantage function is typically computed based on the cumulative return of an entire trajectory. This coarse-grained signal poses a significant credit assignment challenge, as

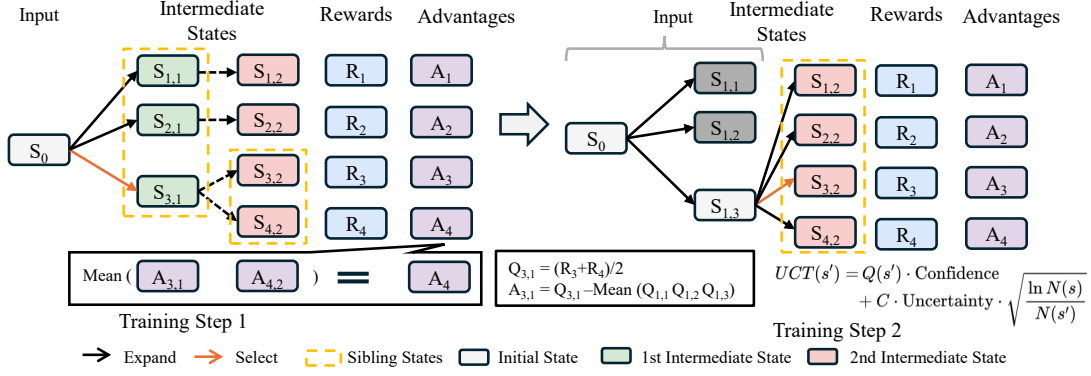


Figure 3: Overview of the M-GRPO algorithm. Starting from the initial state S_0 , the framework expands four simulation sequences in parallel. The **advantage values** of nodes within these sequences are calculated to guide the model training step. Post-training, the **pivotal intermediate waypoint** is selected based on the reformulated UCT formula, serving as the root for the subsequent expansion and training cycle. This iterative process continues until the MCTS termination criterion is satisfied.

it fails to disambiguate the individual contributions of intermediate states. Consequently, it is difficult for the model to pinpoint which specific choices are most critical for achieving success.

To overcome this limitation, we introduce a fine-grained advantage function calculated at the intermediate state level. Our approach is tailored for a tree-search process wherein a LLM generates a set of M candidate sequences (or completions), $\{\tau_1, \dots, \tau_M\}$, for a given planning problem. Each trajectory τ_m is composed of a sequence of intermediate states, $\tau_m = (s_{m,1}, s_{m,2}, \dots, s_{m,T_m})$.

Let $s_{m,n}$ be the n -th intermediate state in the m -th generated sequence. We estimate its corresponding state-value, or Q-value $Q_{m,n}$, as the mean empirical return from all Monte Carlo simulations that traverse this state. Specifically, if $W_{m,n}$ is the sum of cumulative rewards from all visits to state $s_{m,n}$ and $N_{m,n}$ is its total visit count, the Q-value is given by:

$$Q_{m,n} = \frac{W_{m,n}}{N_{m,n}} \quad (5)$$

We then define the advantage of a specific state, $A_{m,n}$, relative to its "sibling" states—i.e., the set of other candidate states $\{s_{j,n}\}_{j=1}^M$ that share a common prefix sequence. The state-level advantage is formulated as:

$$A_{m,n} = Q_{m,n} - \text{Mean}(Q_{\text{siblings}}) \quad (6)$$

where the second term represents the mean Q-value across all sibling states at depth n . This formulation directly quantifies how much better the choice leading to $s_{m,n}$ is compared to the average of alternative choices at that decision point. A deliberate

design choice is the omission of reward normalization (i.e., skipping division by the standard deviation). As the LLM often generates identical optimal completions, forgoing normalization prevents "reward hacking," where the value of a superior path could be artificially deflated due to its high frequency of generation. The necessity of this non-normalization approach in such scenarios has been formally demonstrated in (Liu et al., 2025).

Finally, to align with the GRPO framework, we compute a single advantage value for each trajectory by averaging the advantages of all its constituent intermediate states. For a trajectory τ_m of length T_m , its overall advantage $A(\tau_m)$ is calculated as:

$$A(\tau_m) = \frac{1}{T_m} \sum_{n=1}^{T_m} A_{m,n} \quad (7)$$

This trajectory-level advantage $A(\tau_m)$ is then used as the training signal within the GRPO loss function. This fine-grained approach to advantage calculation provides a more precise and informative signal, enabling the model to learn not only which overall sequences are effective, but also to discern the value of the specific intermediate states that are most critical for constructing an optimal plan. We present the full pseudo-code in Algorithm 1

3 Experiments

3.1 Experimental Setup

Given that spatial reasoning in embodied intelligence often operates under stringent computational constraints, we focus on open-source models with

manageable parameter scales. Specifically, we employ Qwen3-4B-Instruct-2507 as our base model, which is further optimized via our M-GRPO framework to serve as the high-level planner. The untrained version of the same model served as both the environment decomposition model and the action generation model. We deliberately exclude large-scale proprietary closed-source models as they are often unsuitable for real-time embodied tasks, where low-latency local execution, robust offline functionality, and strict data privacy are paramount.

Datasets We evaluate our HSRL framework across four planning benchmarks with increasing difficulty to assess its performance and generalization capability. First, we use Maze Navigation(Valmeekam et al., 2023), a classical task consisting of 1,090 10×10 grids. To create highly challenging scenarios, we configured 40% of the cells in each map as obstacles. The dataset is partitioned into 668 training and 422 testing instances. To validate the real-world applicability of our approach, inspired by (Zhao et al., 2025), we leverage the R2V dataset (Liu et al., 2017), which comprises 815 authentic architectural floorplans. We randomly selected a subset of 50 maps from this collection. Each floorplan was converted into a textual representation and scaled proportionally to a 20×20 resolution, followed by the random sampling of three pairs of start and goal positions for each layout. This procedure resulted in the evaluation benchmark for real-world indoor navigation tasks. Second, to assess out-of-distribution (OOD) generalization, we employ the Blocksworld benchmark(Saha et al., 2025), whose test set is intentionally more complex, featuring more blocks and requiring longer plans (7–10 steps) than the training set (1–6 steps). Finally, we validate our framework on the novel and highly challenging GameTraversalBenchmark (GTB)(Nasir et al., 2024). This benchmark contains 150 diverse maps with multiple objectives and paths exceeding 100 steps. As GTB lacks a training set, we evaluate our Maze-trained model in a zero-shot transfer setting to test its capabilities on complex, unseen tasks.

Baselines We compare HSRL against a diverse set of representative baselines. First, we compare it with foundational reasoning strategies, including the classic Chain-of-Thought (CoT)(Wei et al., 2022a) and ReAct(Yao et al., 2023c), which interleaves reasoning traces with actions for im-

proved synergy. We also include advanced reasoning and self-reflection methods like Inner Monologue(Huang et al., 2023), which enhances internal thought processes, and Reflexion(Shinn et al., 2023), which uses iterative self-correction to refine plans. For direct planning, we use Prog-Prompt(Singh et al., 2023) as a strong representative of in-context learning-based approaches. Furthermore, we contrast HSRL with search-based methods like Tree Planner(Hu et al., 2024) and the hierarchical planner HyperTree(Gui et al.), the latter of which is known to have limitations on spatial reasoning tasks. Finally, we include System-1.x(Saha et al., 2025), a powerful baseline meticulously fine-tuned on tasks similar to ours, which employs a controller to switch between "fast-thinking" and "slow-thinking" modes.

Evaluation metrics We evaluate our model’s planning ability using metrics tailored to each benchmark. For the classical Maze and Blocksworld tasks, we measure the Completion Rate (CR), which is the percentage of successfully solved instances, and the Optimal Rate (OR), defined as the percentage of completed tasks where the plan length matches the shortest path computed by an A* search. For the more complex Game-TraversalBenchmark (GTB), we adopt its official metrics. The primary metric is the GTB Score, a composite measure that assesses performance based on goal proximity, path length, and generation errors (see Appendix B.1). Additionally, we report Top-5 Accuracy, the fraction of tasks where the agent ends within five tiles of the target, to evaluate success in large-scale maps.

3.2 Implementation Detail

M-GRPO Finetuning. We fine-tuned the model using the trl library on eight NVIDIA A800 GPUs. The optimization was conducted using the AdamW optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.999$. We employed a learning rate of 1×10^{-6} following a cosine decay schedule. The hyperparameters τ is 1.0 and γ is 0.4. The Epoch Number was 1 and the Batch Size was 1. For the inference phase, the Temperature was set to 1.0 and the Num generations was 8.

Reward Function. To ensure training stability and prevent reward hacking, we implement a comprehensive reward function that weights task success against path length through multiple shaping terms. Specifically, to mitigate the tendency of RL

agents to over-generate—a common phenomenon where models produce excessively long responses to maximize cumulative rewards—we introduce a stringent length penalty to encourage concise and effective reasoning. Detailed formulations of these composite reward functions are provided in Appendix C.1.

3.3 Results Analysis

The experiments results, detailed in Table 1, clearly demonstrate the effectiveness of our hierarchical planning and search framework.

HSRL significantly improves task performance.

Our method HSRL achieves state-of-the-art (SOTA) performance across all metrics on all tasks. For instance, in the Maze (10 × 10) task, our method achieves a goal completion rate of 61.37%, markedly outperforming other methods such as System-1.x (54.74%) and ReAct (53.80%). On the real-world R2V dataset, given that R2V has a lower obstacle density than the synthetic Maze task, HSRL demonstrates even stronger performance, reaching a completion rate of 62.67%. The larger spatial scale of R2V task typically poses a significant challenge for conventional methods. Specifically, as the search space expands, the reasoning capabilities of monolithic methods tend to deteriorate rapidly. In contrast, HSRL maintains high efficacy by decomposing long-horizon navigation into manageable sub-tasks, ensuring that each segment remains within the model’s optimal problem-solving range. This advantage is even more pronounced in the more complex Blocksworld task, where our method’s 30.50% completion rate far exceeds all baselines. This pattern of superiority extends to the challenging GameTraversalBenchmark (GTB), where HSRL achieves the highest GTB Score of 32.69, surpassing strong competitors like Reflexion (29.34). Furthermore, it also secures the best Top 5 Accuracy at 44.41%, demonstrating its robustness in complex, large-scale planning environments.

Superior Solution Optimality. Beyond merely completing a task, the ability to generate optimal (or near-optimal) paths is a crucial measure of a planning model’s intelligence. On the optimal rate metric for the Maze and Blocksworld tasks, our method achieves optimality rates of 47.87% and 55.33%, respectively, the highest among all compared methods. This result indicates that by combining the forward-search capabilities of MCTS with the general knowledge of large models, our

framework can explore the solution space more thoroughly, effectively avoiding local optimal to devise more efficient and concise solutions.

Cross-Task Robustness and Generalization.

The value of a general-purpose planning model lies in its cross-task generalization capability. As shown in the table, our method performs exceptionally well on the classical spatial reasoning tasks of Maze and Blocksworld, and the complex world knowledge required for the Game Travel Benchmark. In contrast, some baselines exhibit strong task-specific biases; for example, while ReAct performs reasonably well in Maze, its completion rate plummets to just 3.00% in Blocksworld. This comparison validates the robustness of our hierarchical framework, which consistently decomposes complex problems into manageable sub-goals for effective problem-solving, irrespective of the task modality. Moreover, our excellent performance in Blocksworld demonstrates strong out-of-distribution generalization capabilities.

3.4 Further Analysis

Ablation Study. To validate the contribution of each component, we conducted a comprehensive ablation study (Table 2). The results reveal a clear hierarchy of importance. Removing the MCTS module—thereby reverting M-GRPO to a standard GRPO training paradigm (HSRL w/o MCTS)—leads to a notable decline in both success and optimality, confirming that its systematic, forward-looking search is crucial for exploring diverse solution pathways. Further removing the M-GRPO policy optimization, where the base model is used directly as the high-level planner without further training (HSRL Untrained) causes a precipitous performance collapse, especially in optimality (e.g., Maze optimality plummets from 47.87% to 45.97%). This demonstrates that M-GRPO is the core engine that translates the rich search experience from MCTS into a refined planning intuition, endowing the LLM with the ability to generate high-quality, task-aligned sub-goals. Finally, the performance of the State-Hierarchical Only configuration—which eliminates environment-level hierarchy and relies on the base model for both intermediate waypoint identification and low-level action generation—still significantly surpasses direct answering methods, and the inclusion of environment-hierarchical approaches effectively improves task completion.

Table 1: Main results: Comparison of HSRL against baseline methods across various spatial reasoning benchmarks. **CR** and **OR** denote Completion Rate and Optimality Rate, respectively. The best performance in each category is highlighted in **bold**.

Method	Maze (10×10)		R2V (Real-World)		Blocksworld (5-7)		GTB	
	CR (%) ↑	OR (%) ↑	CR (%) ↑	OR (%) ↑	CR (%) ↑	OR (%) ↑	Score ↑	Top-5 Acc. ↑
Direct Answer	23.69	23.45	1.33	1.13	6.50	6.00	23.61	25.96
CoT	43.12	38.39	16.67	12.33	10.50	8.00	26.58	31.61
Reflexion	45.02	37.91	22.00	21.33	15.00	8.50	29.34	37.76
ReAct	53.80	26.06	16.00	13.33	8.00	3.00	20.40	39.85
ProgPrompt	34.60	33.41	43.33	42.00	9.50	6.00	22.45	26.12
Inner Monologue	54.03	34.60	27.33	16.67	4.00	0.00	19.18	21.41
System-1.x	54.74	36.02	50.67	49.33	27.00	14.50	27.73	30.01
HyperTree	37.91	22.98	25.33	19.67	8.00	3.50	25.81	26.67
Tree Planner	39.10	27.01	6.33	3.00	7.00	4.00	25.28	25.46
HSRL (Ours)	61.37	47.87	62.67	55.33	30.50	21.00	32.69	44.41

Table 2: Ablation study on the core components of our HSRL framework. We evaluate the contribution of each module across various benchmarks. **CR** and **OR** represent Completion Rate and Optimality Rate, respectively. The best performance is highlighted in **bold**.

Model Configuration	Maze (10×10)		R2V (Real-World)		Blocksworld (5-7)		GTB	
	CR (%) ↑	OR (%) ↑	CR (%) ↑	OR (%) ↑	CR (%) ↑	OR (%) ↑	Score ↑	Top-5 Acc. ↑
State-Hierarchical Only	50.24	13.03	52.33	25.67	10.00	9.50	26.16	29.85
HSRL (Untrained)	54.50	14.22	60.33	28.00	12.50	9.50	26.96	32.38
HSRL (w/o MCTS)	55.21	45.97	58.00	54.67	28.00	15.00	27.80	38.09
HSRL (Ours)	61.37	47.87	62.67	55.33	30.50	21.00	32.69	44.41

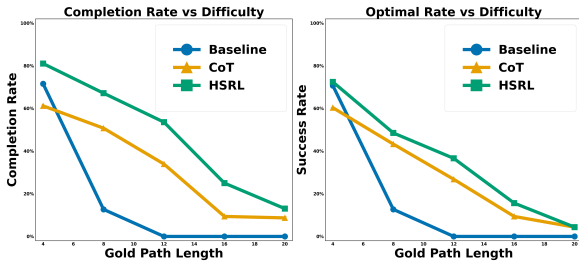


Figure 4: Performance degradation as task difficulty increases. HSRL shows greater robustness.

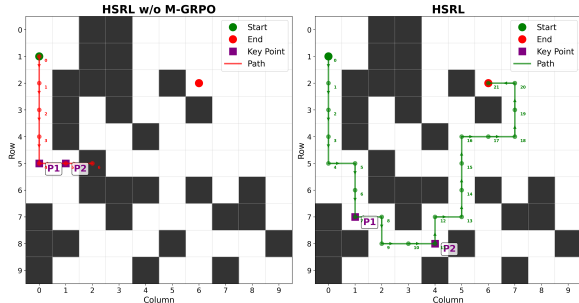


Figure 5: Qualitative comparison of a generated plan before (left) and after (right) M-GRPO training.

Robustness and Qualitative Insights. Further analysis highlights HSRL’s robustness. As shown in Figure 4, HSRL’s performance degrades far more gracefully with increasing task difficulty compared

to baselines, maintaining a substantial and reliable advantage on the most challenging Maze instances. This resilience is not merely statistical; it stems from a fundamental improvement in high-level planning quality. A qualitative comparison in Figure 5 illustrates the underlying mechanism: while an untrained model generates ill-conceived sub-goals leading to a failed plan, the M-GRPO trained HSRL produces a strategic and successful path. This synergy between quantitative robustness and qualitative intelligence validates our framework’s effectiveness in complex planning scenarios where long-horizon reasoning is paramount.

4 Conclusion

This paper introduced HSRL, a hierarchical framework that overcomes LLM spatial reasoning limitations through state and environment decomposition. By combining MCTS-guided GRPO algorithm, HSRL achieves SOTA performance across navigation and strategic benchmarks. Our results show substantial gains in completion rates and path optimality over existing methods. This work establishes a scalable approach for integrating LLMs into complex physical environments and embodied intelligence tasks.

570 Limitations

571 Despite the advancements in spatial reasoning
572 achieved by this study, several limitations remain.
573 First, M-GRPO introduces additional computa-
574 tional overhead during the training phase. Since
575 integrating MCTS for policy optimization necessi-
576 tates multiple path simulations, the training process
577 is more time-consuming than conventional fine-
578 tuning. However, this overhead is confined to the
579 offline stage and does not compromise inference
580 efficiency. Second, due to the scarcity of large-
581 scale datasets featuring pure textual descriptions of
582 3D environments, our experiments currently focus
583 on 2D indoor navigation tasks (R2V). Future work
584 will explore the generalization of this framework
585 to more complex 3D spatial representations and
586 multimodal environments.

587 References

- 588 Mohamed Aghzal, Erion Plaku, and Ziyu Yao. 2024a.
589 [Can large language models be good path planners? a
590 benchmark and investigation on spatial-temporal rea-
591 soning](#). In *ICLR 2024 Workshop on Large Language
592 Model (LLM) Agents*.
- 593 Mohamed Aghzal, Erion Plaku, and Ziyu Yao. 2024b.
594 [Look further ahead: Testing the limits of gpt-4 in
595 path planning](#). In *2024 IEEE 20th International
596 Conference on Automation Science and Engineering
597 (CASE)*, pages 1020–1027. IEEE.
- 598 Amos Azaria and Tom Mitchell. 2023. [The internal
599 state of an llm knows when it’s lying](#). *arXiv preprint
600 arXiv:2304.13734*.
- 601 Yongchao Chen, Jacob Arkin, Charles Dawson, Yang
602 Zhang, Nicholas Roy, and Chuchu Fan. 2024. [Auto-
603 tamp: Autoregressive task and motion planning with
604 llms as translators and checkers](#). In *2024 IEEE In-
605 ternational conference on robotics and automation
606 (ICRA)*, pages 6695–6702. IEEE.
- 607 Alan Dao and Dinh Bach Vu. 2025. [Alphamaze: En-
608 hancing large language models’ spatial intelligence
609 via grpo](#). *arXiv preprint arXiv:2502.14669*.
- 610 Hourui Deng, Hongjie Zhang, Jie Ou, and Chaosheng
611 Feng. 2025. [Can llm be a good path planner based
612 on prompt engineering? mitigating the hallucination
613 for path planning](#). In *International Conference on
614 Intelligent Computing*, pages 3–15. Springer.
- 615 J. Duan, S. E. Li, Y. Guan, Q. Sun, and B. Cheng. 2020.
616 [Hierarchical reinforcement learning for self-driving
617 decision-making without reliance on labelled driving
618 data](#). *IET Intelligent Transport Systems*, 14:297–305.

- Lutfi Eren Erdogan, Hiroki Furuta, Sehoon Kim,
Nicholas Lee, Suhong Moon, Gopala Anu-
manchipalli, Kurt Keutzer, and Amir Gholami. 2025a.
[Plan-and-act: Improving planning of agents for long-
horizon tasks](#). In *Forty-second International Confer-
ence on Machine Learning*. 619
620
621
622
623
624
- Lutfi Eren Erdogan, Nicholas Lee, Sehoon Kim, Suhong
Moon, Hiroki Furuta, Gopala Anumanchipalli, Kurt
Keutzer, and Amir Gholami. 2025b. [Plan-and-act:
Improving planning of agents for long-horizon tasks](#).
arXiv preprint arXiv:2503.09572. 625
626
627
628
629
- Nanxu Gong, Chandan K. Reddy, Wangyang Ying,
Haifeng Chen, and Yanjie Fu. 2024. [Evolutionary
large language model for automated feature transfor-
mation](#). *Preprint*, arXiv:2405.16203. 630
631
632
633
- Xinyu Guan, Li Lyna Zhang, Yifei Liu, Ning Shang,
Youran Sun, Yi Zhu, Fan Yang, and Mao Yang. [rstar-
math: Small llms can master math reasoning with
self-evolved deep thinking](#). In *Forty-second Interna-
tional Conference on Machine Learning*. 634
635
636
637
638
- Runquan Gui, Zhihai Wang, Jie Wang, Chi Ma, Huiling
Zhen, Mingxuan Yuan, Jianye HAO, Defu Lian, En-
hong Chen, and Feng Wu. [Hypertree planning: En-
hancing llm reasoning via hierarchical thinking](#). In
*Forty-second International Conference on Machine
Learning*. 639
640
641
642
643
644
- Runquan Gui, Zhihai Wang, Jie Wang, Chi Ma, Huil-
ing Zhen, Mingxuan Yuan, Jianye HAO, Defu Lian,
Enhong Chen, and Feng Wu. 2025. [Hypertree plan-
ning: Enhancing LLM reasoning via hierarchical
thinking](#). In *Forty-second International Conference
on Machine Learning*. 645
646
647
648
649
650
- Mengkang Hu, Yao Mu, Xinmiao Chelsey Yu, Mingyu
Ding, Shiguang Wu, Wenqi Shao, Qiguang Chen, Bin
Wang, Yu Qiao, and Ping Luo. 2024. [Tree-planner:
Efficient close-loop task planning with large language
models](#). In *The Twelfth International Conference on
Learning Representations*. 651
652
653
654
655
656
- Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky
Liang, Pete Florence, Andy Zeng, Jonathan Tompson,
Igor Mordatch, Yevgen Chebotar, and 1 others. 2023.
[Inner monologue: Embodied reasoning through plan-
ning with language models](#). In *Conference on Robot
Learning*, pages 1769–1782. PMLR. 657
658
659
660
661
662
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023.
[Semantic uncertainty: Linguistic invariances for un-
certainty estimation in natural language generation](#).
In *The Eleventh International Conference on Learn-
ing Representations*. 663
664
665
666
667
- Harsh Kumar, Ruiwei Xiao, Benjamin Lawson, Ilya
Musabirov, Jiakai Shi, Xinyuan Wang, Huayin Luo,
Joseph Jay Williams, Anna N. Rafferty, John Stamper,
and Michael Liut. 2024. [Supporting self-reflection
at scale with large language models: Insights from
randomized field experiments in classrooms](#). In *Pro-
ceedings of the Eleventh ACM Conference on Learn-
ing @ Scale, L@S ’24*, page 86–97, New York, NY,
USA. Association for Computing Machinery. 668
669
670
671
672
673
674
675
676

677	Chen Liu, Jiajun Wu, Pushmeet Kohli, and Yasutaka Furukawa. 2017. Raster-to-vector: Revisiting floor-plan transformation. In <i>Proceedings of the IEEE international conference on computer vision</i> , pages 2195–2203.	733
678		734
679		735
680		
681		
682	Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. 2025. Understanding r1-zero-like training: A critical perspective. <i>Preprint</i> , arXiv:2503.20783.	736
683		737
684		738
685		739
686	Xinyang Lu, Flint Xiaofeng Fan, and Tianying Wang. 2023. Action and trajectory planning for urban autonomous driving with hierarchical reinforcement learning. In <i>ICML Workshop on New Frontiers in Learning, Control, and Dynamical Systems</i> .	740
687		741
688		742
689		743
690		744
691	Runyu Ma, Jelle Luijkx, Zlatan Ajanović, and Jens Kober. 2025. Explorllm: Guiding exploration in reinforcement learning with large language models. In <i>2025 IEEE International Conference on Robotics and Automation (ICRA)</i> , pages 9011–9017. IEEE.	745
692		746
693		747
694		748
695		
696	Suvir Mirchandani, Fei Xia, Pete Florence, Brian Ichter, Danny Driess, Montserrat Gonzalez Arenas, Kanishka Rao, Dorsa Sadigh, and Andy Zeng. 2023. Large language models as general pattern machines. In <i>Proceedings of The 7th Conference on Robot Learning</i> , volume 229 of <i>Proceedings of Machine Learning Research</i> , pages 2498–2518. PMLR.	749
697		750
698		751
699		752
700		753
701		754
702		
703	Ida Momennejad, Hosein Hasanbeig, Felipe Vieira Frujeri, Hiteshi Sharma, Nebojsa Jojic, Hamid Palangi, Robert Ness, and Jonathan Larson. 2023. Evaluating cognitive maps and planning in large language models with cogeval. In <i>Advances in Neural Information Processing Systems</i> , volume 36, pages 69736–69751. Curran Associates, Inc.	755
704		756
705		757
706		758
707		759
708		760
709		
710	Muhammad Umair Nasir, Steven James, and Julian Togelius. 2024. Gametraversalbenchmark: Evaluating planning abilities of large language models through traversing 2d game maps. In <i>Advances in Neural Information Processing Systems</i> , volume 37, pages 31813–31827. Curran Associates, Inc.	761
711		762
712		763
713		764
714		765
715		766
716	Matthew Renze and Erhan Guven. 2024. Self-reflection in llm agents: Effects on problem-solving performance. <i>CoRR</i> , abs/2405.06682.	767
717		768
718		769
719	Swarnadeep Saha, Archiki Prasad, Justin Chen, Peter Hase, Elias Stengel-Eskin, and Mohit Bansal. 2025. System 1.x: Learning to balance fast and slow planning with language models. In <i>The Thirteenth International Conference on Learning Representations</i> .	770
720		771
721		772
722		773
723		774
724	Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik R Narasimhan, and Shunyu Yao. 2023. Reflexion: language agents with verbal reinforcement learning. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	775
725		776
726		777
727		778
728		779
729	Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. 2023. Prog-prompt: Generating situated robot task plans using large language models. In <i>2023 IEEE International Conference on Robotics and Automation (ICRA)</i> , pages 11523–11530. IEEE.	780
730		781
731		782
732		783
	Karthik Valmeekam, Matthew Marquez, Sarath Sreedharan, and Subbarao Kambhampati. 2023. On the planning abilities of large language models - a critical investigation. In <i>Advances in Neural Information Processing Systems</i> , volume 36, pages 75993–76005. Curran Associates, Inc.	784
		785
		786
		787
		788
		789
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022a. Chain-of-thought prompting elicits reasoning in large language models. In <i>Advances in Neural Information Processing Systems</i> , volume 35, pages 24824–24837. Curran Associates, Inc.	790
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800
		801
		802
		803
		804
		805
		806
		807
		808
		809
		810
		811
		812
		813
		814
		815
		816
		817
		818
		819
		820
		821
		822
		823
		824
		825
		826
		827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

790	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak	method (Wu et al., 2024) materializes intermediate	841
791	Shafraan, Karthik R Narasimhan, and Yuan Cao.	states to assist with reasoning. Nevertheless, this	842
792	2023c. React: Synergizing reasoning and acting	iterative feedback loop often results in high costs	843
793	in language models . In <i>The Eleventh International</i>	and inefficiency in querying or interactions.	844
794	<i>Conference on Learning Representations</i> .		
795	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,	A.2 Hierarchical Method	845
796	Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen	Hierarchical reasoning breaks down decision-	846
797	Zhang, Junjie Zhang, Zican Dong, and 1 others. 2023.	making tasks into multiple levels, from high-level	847
798	A survey of large language models. <i>arXiv preprint</i>	strategic planning to low-level specific control.	848
799	<i>arXiv:2303.18223</i> .	This decomposition reduces computational com-	849
800	Yubo Zhao, Qi Wu, Yifan Wang, Yu-Wing Tai, and Chi-	plexity by solving several less difficult sub-tasks,	850
801	Keung Tang. 2025. LLM-navi: Navigating motion	thus enabling the handling of tasks more challeng-	851
802	agents in cluttered and dynamic environments via	ing than direct complex reasoning. Hierarchical	852
803	LLM reasoning .	reasoning has achieved notable results in many re-	853
804	Wei Zhu and Mitsuhiro Hayashibe. 2023. A hierarchi-	inforcement learning tasks, especially in embodied	854
805	cal deep reinforcement learning framework with high	AI scenarios. For example, (Duan et al., 2020)	855
806	efficiency and generalization for fast and safe navi-	has applied hierarchical methods to autonomous	856
807	gation . <i>IEEE Transactions on Industrial Electronics</i> ,	driving, allowing for smooth and safe decision-	857
808	70(5):4962–4971.	making on highways. (Lu et al., 2023) and (Zhu	858
809	Appendix	and Hayashibe, 2023) separate decision-making	859
810	A Related Work	tasks into different layers, such as global path plan-	860
811	A.1 Spatial Reasoning in large language	ning and local motion control. These models ben-	861
812	models	efit from breaking down the decision-making pro-	862
813	Many researchers have pointed out that LLMs have	cess into simpler, more tractable components, en-	863
814	weaknesses in spatial reasoning or spatial plan-	abling each layer to focus on a specific task. This	864
815	ning(Aghzal et al., 2024a,b). To address these	enhances computational efficiency and decision ac-	865
816	issues, some methods leverage in-context exam-	curacy in complex environments.	866
817	ples and prompting techniques, such as Chain-	In recent years, hierarchical reasoning methods	867
818	of-Thought (CoT)(Wei et al., 2022a) and Tree-	have also been successfully introduced into the	868
819	of-Thought (ToT)(Yao et al., 2023b), which have	planning tasks of LLMs. For instance, DeAR(Xue	869
820	demonstrated remarkable reasoning abilities in var-	et al., 2024) imitates the human reasoning cycle	870
821	ious tasks. However, for spatial reasoning tasks,	by using a tree-based question decomposition ap-	871
822	in-context learning often fails because LLMs lack	proach to organize the reasoning process and break	872
823	spatial reasoning knowledge or their knowledge	down problems into simpler sub-questions. Hyper-	873
824	even conflicts with it.	Tree Planning(Gui et al., 2025) is a new paradigm	874
825	To overcome this challenge, some studies utilize	that enhances LLM reasoning with a hypertree	875
826	LLMs for general-purpose reasoning, converting	structure. It effectively breaks down intricate rea-	876
827	spatial information into logical forms(Yang et al.,	soning steps using a flexible divide-and-conquer	877
828	2023) or using them as a general pattern machine	strategy to handle diverse constraints and manage	878
829	for sequence transformation(Mirchandani et al.,	multiple distinct sub-tasks, demonstrating superior	879
830	2023; Gong et al., 2024). Recently, other works	performance in complex tasks like travel planning.	880
831	have evaluated LLMs as a cognitive capability in	Plan-and-Act(Erdogan et al., 2025a) explicitly se-	881
832	navigation and planning tasks(Mommenjad et al.,	parates high-level planning from low-level execu-	882
833	2023). However, these methods perform poorly in	This framework includes a PLANNER model for	883
834	tasks requiring continuous action reasoning.	generating structured high-level plans and an EX-	884
835	Another mainstream approach introduces closed-	ECUTOR model for translating these plans into	885
836	loop feedback mechanisms. Some works, like	environment-specific actions, thereby improving	886
837	(Renze and Guven, 2024), use self-reflection for	performance on complex multi-step tasks such as	887
838	self-evaluation and replanning, while others adopt	web navigation.	888
839	external feedback for reflection (Kumar et al.,	However, these methods only consider high-	889
840	2024). Furthermore, the Vision-of-Thought (VoT)	level, coarse-grained task planning and do not fully	890
		leverage the potential of hierarchical reasoning for	891

low-level tasks that require fine motion control, such as robotic arm motion planning. Therefore, this study aims to fill this gap by solving the complex action planning problem.

B More Background

B.1 GTB Score

$$\text{GTB_Score} = \frac{1}{M} \sum_{m=1}^M \frac{1}{\Delta R^{(m)}} \left(R^{(m)} - \text{LLM}_{PL}^{(m)} - \varepsilon^{(m)} - R_{\min}^{(m)} \right) \quad (8)$$

where:

- M = total number of maps in the dataset.
- $R^{(m)}$ = reward obtained for map m , determined by the final distance d to the objective:

$$R^{(m)} = \begin{cases} +200, & d = 0 \\ +100, & d = 1 \\ +50, & d \in [2, 3] \\ +25, & d \in [3, 5] \\ -50, & d \in [5, 8] \\ -100, & d \geq 8 \end{cases}$$

- $\text{LLM}_{PL}^{(m)}$ = path length taken by the LLM agent on map m .
- $\varepsilon^{(m)}$ = total generation errors made by the LLM on map m .
- $R_{\max}^{(m)} = 200 - A_{PL}^*(m)$, the maximum achievable reward (perfect path with no errors), where $A_{PL}^*(m)$ is the optimal path length computed by an A^* agent.
- $R_{\min}^{(m)} = -100 - A_{PL}^*(m) - \varepsilon_{\max}^{(m)}$, the minimum achievable reward (farthest position, maximal path cost, and maximal errors).

C MORE IMPLEMENTATION DETAILS

C.1 M-GRPO Reward

For a sampled completion $completion_i$, we parse its anchor list $A_i = [a_{i1}, a_{i2}, \dots, a_{in}]$. If the anchor list cannot be parsed, we directly assign a fixed penalty; otherwise, the reward score is computed by a signed power transformation to enlarge

the margin between high- and low-quality completions:

$$R_i = \begin{cases} \text{PARSE_FAIL_PENALTY}, & \text{if } A_i = \emptyset, \\ \text{sign}(z_i) |z_i|^p, & \text{otherwise,} \end{cases} \quad (9)$$

where $\text{sign}(z_i)$ preserves the direction of z_i , and $p > 1$ amplifies its magnitude non-linearly.

The raw score z_i aggregates the quality of anchors visited by a trajectory, while discouraging the use of overly many anchors through a penalty term:

$$z_i = \sum_{a \in A_i} \bar{r}(a) - \alpha \cdot \max(0, |A_i| - A_{\text{expected}}). \quad (10)$$

To evaluate each anchor consistently, we first assign each completion an initial reward r_i according to its alignment with the Manhattan distance of the optimal A^* path. :

$$r_i = \text{BASIC_QUALITY_SCORE} - \left| \sum_{a_i \in A_i} |a_{i+1} - a_i| - \sum_{\hat{a}_i \in A^*} |a_{i+1} - \hat{a}_i| \right|. \quad (11)$$

Each anchor reward $\bar{r}(a)$ is then defined as the average quality of all completions that pass through it, reflecting a consensus measure across different trajectories:

$$\bar{r}(a) = \frac{\sum_{i: a \in A_i} r_i}{|\{i \mid a \in A_i\}|}. \quad (12)$$

D Declaration of Use of AI Assistants

During the preparation of this work, the authors used Gemini-3-Flash to polish the manuscript for language clarity and to assist with LaTeX formatting. Following the use of this tool, the authors reviewed and edited the content as needed and take full responsibility for the final version of the paper.

E Prompts and Examples

High-level Planning Prompt for Maze

Role

You are an expert high-level path planner. You must strictly adhere to the requirements outlined in the system instructions and tasks I have provided to you.

Instructions

1. Your task is to plan a feasible, obstacle-free path for a single agent in a given 10x10 grid environment, from a start to an end point.
2. The path should be defined by a series of key anchor point coordinates.
3. You must identify exactly `{num_anchors}` feasible intermediate anchor points for the given task. These anchor points should be the key turning points of the path used to navigate around obstacles or toward the goal.

Anchor Point Selection Strategy

- The path does not have to be the shortest path. The priority is feasibility and safety (avoiding all obstacles).
- Explore multiple valid paths and select a reasonable one to define your anchor points.
- Anchor points should be strategically located at important positions around obstacles.

Output Format

- You must strictly follow the format below to output the list of anchor points.
- Do not provide any explanation or text other than the final trajectory list.
- Directly output the result in the given format:

```
<trajectory for planning> = [(start_x, start_y), (anchor_1_x, anchor_1_y), ..., (end_x, end_y)]
```

Examples

Example 1:

Task: You are in a 10 by 10 world. There are obstacles that you must avoid at: (4,7), (8,6), (3,3), (9,5), (8,9), (1,1), (5,4), (1,3), (9,9), (4,1), (5,7), (1,6), (9,0), (8,3), (0,0), (7,1), (4,6), (5,0), (2,5) and (4,0). Go from (2,1) to (0,2).
<trajectory for planning> = [(2,1),(2,2),(1,2),(0,2)]

Example 2:

Task: You are in a 10 by 10 world. There are obstacles that you must avoid at: (0,7), (3,2), (0,4), (3,4), (4,6), (7,2), (7,3), (2,0), (3,9), (9,3), (8,2), (9,5), (8,4), (7,5), (4,8), (5,2), (5,5), (7,8), (6,3) and (9,8). Go from (6,8) to (6,1).
<trajectory for planning> = [(6,0),(6,4),(5,3),(6,1)]

Example 3:

Task: You are in a 10 by 10 world. There are obstacles that you must avoid at: (8,5), (7,2), (1,7), (2,0), (3,2), (5,0), (1,9), (3,3), (3,6), (4,7), (0,3), (5,7), (5,3), (4,6), (2,8), (4,3), (9,0), (7,5), (5,5) and (8,9). Go from (0,8) to (7,1).
<trajectory for planning> = [(0,8),(0,4),(1,1),(7,1)]

Task to Solve

Task: You are in a 10 by 10 world. There are obstacles that you must avoid at: (9,6), (1,0), (3,7), (4,4), (9,1), (4,0), (3,4), (8,9), (7,1), (5,1), (3,6), (4,9), (4,8), (0,1), (6,4) and (0,0). Go from (4,1) to (6,8).
<trajectory for planning> = [(4,1),(5,3),(5,6),(6,8)]

Environment Decomposition Prompt for Maze

Given the following list of coordinate points:

```
[(9,6), (1,0), (3,7), (4,4), (9,1), (4,0), (3,4), (8,9), (7,1), (5,1), (3,6), (4,9), (4,8), (0,1), (6,4), (0,0)]
```

Please select all points from the list that satisfy the following condition:

The x-coordinate (x) of the point must be within the closed interval from 4 to 5.

Please strictly follow this format for the output, including only the selected points:

Output Format

```
<obstacles> = [ (x1, y1), (x2, y2), ..., (xn, yn)]
```

```
<assistant>
```

```
<obstacles> = [(4,4),(4,0),(5,1),(4,9),(4,8)]
```

Environment Decomposition Prompt for Maze

Given the following list of coordinate points:

```
[(9,6), (1,0), (3,7), (4,4), (9,1), (4,0), (3,4), (8,9), (7,1), (5,1), (3,6), (4,9), (4,8), (0,1), (6,4), (0,0)]
```

Please select all points from the list that satisfy the following condition:

The y-coordinate (y) of the point must be within the closed interval from 1 to 3.

Please strictly follow this format for the output, including only the selected points:

Output Format

```
<obstacles> = [ (x1, y1), (x2, y2), ..., (xn, yn)]
```

```
<assistant>
```

```
<obstacles> = [(9,1),(7,1),(5,1),(0,1)]
```

Low-level Execution Prompt for Maze

Role

You are an low-level path planner located in a 10 by 10 world. You must strictly adhere to the requirements of the tasks I have provided to you.

Environment

Provide a sequence of actions to navigate a world to reach a goal. (0,0) is located in the upper-left corner and (9,9) lies in down-right corner.

Rules

- <left = (0,-1)>
- <right = (0,+1)>
- <up = (-1,0)>
- <down = (+1,0)>

Output Format

Actions = [action_0 action_1 ... action_n]

Here are some examples:

###

Task: You are in a 10 by 10 world. There are obstacles that you must avoid at: (2,1). Go from (0,1) to (3,4).

Actions = [right right right down down down]

###

Task: You are in a 10 by 10 world. There are obstacles that you must avoid at: (1,5) and (1,2). Go from (5,4) to (0,5).

Actions = [up up up up up right]

###

Task: You are in a 10 by 10 world. There are obstacles that you must avoid at: (0,3), (2,5) and (5,2). Go from (4,2) to (0,5)

Actions = [up up up right right up right]

Task to Solve **subtask_1**

Task: You are in a 10 by 10 world. There are obstacles that you must avoid at:(4,4),(4,0),(5,1),(4,9),(4,8),(9,1),(7,1),(5,1),(0,1). Go from (4,1) to (5,3).

Actions = [right right down]

High-level Planning Prompt for Blocksworld

Role

You are a high-level Blocksworld planner. You must strictly adhere to the requirements outlined in the system instructions and tasks I have provided to you.

Instructions

1. Plan a feasible sequence of block configurations from an initial state to a goal state.
2. Define the plan with exactly `{{num_anchors}}` intermediate stack states.
3. Intermediate states should be key subgoals (e.g., clearing a block or forming partial stacks).

Strategy

- The sequence need not be shortest; feasibility and clarity are the priority.
- Choose anchor states that mark meaningful progress toward the goal.

Output Format

- You must strictly follow the format below to output the list of anchor states.
- Do not provide any explanation or text other than the final output list.
- Directly output the result in the given format:

```
Output = [initial_state, anchor_state1, ..., goal_state]
```

Examples

Example 1:

The initial state:

A is on the table. B is on A. B is clear.

The goal is:

B is on the table. A is on B. A is clear.

```
Output = ['A is on the table. B is on A. B is clear.', 'A is on the table. A is clear. B is on the table. B is clear.', 'B is on the table. A is on B. A is clear.']
```

Example 2:

The initial state:

C is on the table. D is on C. D is clear.

The goal is:

C is on D. D is on the table. D is clear.

```
Output = ['C is on the table. D is on C. D is clear.', 'C is on the table. C is clear. D is on the table. D is clear.', 'C is on D. D is on the table. D is clear.']
```

Example 3:

The initial state:

B is on the table. C is on B. A is on C. A is clear.

The goal is:

B is on the table. B is clear.

C is on B. A is on C. A is clear.

```
Output = ['B is on the table. C is on B. A is on C. A is clear.', 'B is on the table. B is clear. C is on the table. C is clear. A is on the table. A is clear.', 'B is on the table. C is on B. C is clear. A is on the table. A is clear.']
```

Task to Solve

The initial state:

B is on the table. B is clear.

D is on the table. C is on D. E is on C. A is on E. A is clear.

The goal is:

C is on the table. E is on C. D is on E. B is on D. A is on B. A is clear.

```
Output = ['The hand is empty. B is on the table. B is clear. D is on the table. C is on the table. C is clear. E is on the table. E is clear. A is on the table. A is clear.', 'C is on the table. E is on C. D is on E. B is on the table. B is clear. A is on the table. A is clear.']
```

Environment Decomposition Prompt for Blocksworld

Given the 2 state of the same group of Blocks:

<1>'B is on the table. B is clear. D is on the table. C is on the table. C is clear. E is on the table. E is clear. A is on the table. A is clear.'

<2>'C is on the table. E is on C. D is on E. B is on the table. B is clear. A is on the table. A is clear.'

Please select all state pairs (including block-block pairs, block-table pairs, and block-clear pairs) whose relative positions remain unchanged.

Please strictly follow this format for the output, including only the selected pairs:

Output Format

<pairs> = [(Object1,Object2),(Object3,Object4),...]

<assistant>

<pairs> = [(B,table),(A,table),(C,table)]

Low-level Execution Prompt for Blocksworld

You are a Blocks World action planner. You will be given an initial state and a goal state. Provide a sequence of actions to move the blocks to reach a goal similarly to the examples below. Do not include any extra text or explanations.

Here are some examples:

###

The initial state:

The hand is empty.

B is on the table. A is on B. C is on A. C is clear.

The goal is:

B is on the table. A is on B. A is clear.

C is on the table. C is clear.

<Observation>: B is still on table; A is still on B

Actions: Move C from A to table

###

The initial state:

The hand is empty.

B is on the table. C is on B. D is on C. A is on D. A is clear.

The goal is:

A is on the table. A is clear.

C is on the table. B is on C. D is on B. D is clear.

<Observation>: A is still clear

Actions: Move A from D to table | Move D from C to table | Move C from B to table | Move B from table to C | Move D from table to B

Now, here is your task:

###

The initial state:

B is on the table. B is clear.

D is on the table. D is clear.

C is on the table. C is clear.

E is on the table. E is clear.

A is on the table. A is clear.

The goal is:

C is on the table. E is on C. D is on E. D is clear.

B is on the table. B is clear.

A is on the table. A is clear.

<Observation>: B is still on table; A is still on table; C is still on table

Actions: Move E from table to C | Move D from table to E

High-level Planning Prompt for GTB

Role

You are an expert high-level planner for a 2D-grid game. You must strictly adhere to the requirements outlined in the system instructions and tasks I have provided to you.

Instructions

1. Your task is to plan a feasible sequence of moves in a $\{\{world_height\}\} \times \{\{world_length\}\}$ grid environment, from a start state to a goal state.
2. The plan should be defined by a series of key anchor coordinates.
3. You must identify exactly $\{\{num_anchors\}\}$ feasible anchor states for the given task. These should be important turning points or subgoals used to navigate around obstacles or toward objectives.

Anchor State Selection Strategy

- The plan does not need to be the shortest; the priority is feasibility and safety (avoiding all obstacles).
- Anchor states should be strategically located at key subgoals: clearing an obstacle, moving around blocking tiles, or forming partial progress toward objectives.
- Explore multiple valid strategies and select one reasonable plan.

Reward Context

- Rewards are given as follows:
 - $\{\{reward_design\}\}$
 - $\{\{reward_feedback\}\}$
- You are also given information about your previous attempt:
 - Actions generated: $\{\{total_actions[list(objective_tile_dict.keys())[i]]\}$
 - Start position: $\{\{prev_protagonist_position\}\}$
 - End position: $\{\{protagonist_position\}\}$
 - Distance from objective: $\{\{distance_from_objective\}\}$
 - Objective location: $\{\{list(objective_tile_dict.values())[i]\}\}$
 - GTB Reward received: $\{\{reward_this_objective[list(objective_tile_dict.keys())[i]]\}\}$

Output Format

- Strictly output the result in the following format, without any explanation:
<trajectory for planning> = [(start_x, start_y), (anchor_1_x, anchor_1_y), ..., (end_x, end_y)]

Examples

Example 1:

Task: You are in a 10 by 10 world. There are obstacles that you have to avoid at: (4,7), (8,6), (3,3), (9,5), (8,9), (1,1), (5,4), (1,3), (9,9), (4,1), (5,7), (1,6), (9,0), (8,3), (0,0), (7,1), (4,6), (5,0), (2,5) and (4,0). Go from (2,1) to (0,2).
<trajectory for planning> = [(2,1),(2,2),(0,2)]

Example 2:

Task: You are in a 10 by 10 world. There are obstacles that you have to avoid at: (0,7), (3,2), (0,4), (3,4), (4,6), (7,2), (7,3), (2,0), (3,9), (9,3), (8,2), (9,5), (8,4), (7,5), (4,8), (5,2), (5,5), (7,8), (6,3) and (9,8). Go from (6,8) to (6,1).
<trajectory for planning> = [(6,0),(6,4),(5,3),(4,2),(6,1)]

Example 3:

Task: You are in a 10 by 10 world. There are obstacles that you have to avoid at: (8,5), (7,2), (1,7), (2,0), (3,2), (5,0), (1,9), (3,3), (3,6), (4,7), (0,3), (5,7), (5,3), (4,6), (2,8), (4,3), (9,0), (7,5), (5,5) and (8,9). Go from (0,8) to (7,1).
<trajectory for planning> = [(0,8),(0,4),(1,1),(7,1)]

Task to Solve

Task: $\{\{task\}\}$
<trajectory for planning> =

Environment Decomposition Prompt for GTB

Given the following list of coordinate points:
[`{obstacles_this_object}`]

Please select all points from the list that satisfy the following condition:
The x-coordinate (x) of the point must be within the closed interval from `{x_this_object_min}` to `{x_this_object_max}`.

Please strictly follow this format for the output, including only the selected points:
Output Format
`<obstacles> = [(x1, y1), (x2, y2), ..., (xn, yn)]`

`<assistant>`
`<obstacles> =`

Environment Decomposition Prompt for GTB

Given the following list of coordinate points:
[`{obstacles_this_object}`]

Please select all points from the list that satisfy the following condition:
The y-coordinate (y) of the point must be within the closed interval from `{y_this_object_min}` to `{y_this_object_max}`.

Please strictly follow this format for the output, including only the selected points:
`<obstacles> = [(x1, y1), (x2, y2), ..., (xn, yn)]`

`<assistant>`
`<obstacles> =`

Low-level Execution Prompt for GTB

Role

You are a low-level path planner located in a {world_height} by {world_width} world. You must strictly adhere to the requirements of the tasks I have provided to you.

Environment

Provide a sequence of actions to navigate a world to reach a goal. (0,0) is located in the upper-left corner and {(world_height,world_width)} lies in down-right corner.

Rules

- <left = (0,-1)>
- <right = (0,+1)>
- <up = (-1,0)>
- <down = (+1,0)>

Output Format

Actions = [action_0 action_1 ... action_n]

Here are some examples:

###

Task: You are in a 10 by 10 world. There are obstacles that you must avoid at: (2,1). Go from (0,1) to (3,4).

Actions = [right right right down down down]

###

Task: You are in a 10 by 10 world. There are obstacles that you must avoid at: (1,5) and (1,2). Go from (5,4) to (0,5).

Actions = [up up up up up right]

###

Task: You are in a 10 by 10 world. There are obstacles that you must avoid at: (0,3), (2,5) and (5,2). Go from (4,2) to (0,5)

Actions = [up up up right right up right]

Task to Solve **subtask_1**

Task: {GTB_sub_task}

Actions =