

Collision- and Reachability-Aware Multi-Robot Control with Grounded LLM Planners

Jiabao Ji^{1*} Yongchao Chen^{2,3} Yang Zhang⁴
 Ramana Rao Kompella⁵ Chuchu Fan² Gaowen Liu⁵ Shiyu Chang^{1*}

¹UC Santa Barbara ²MIT ³Harvard University ⁴MIT-IBM Watson AI Lab ⁵Cisco Research

Abstract

Large language models (LLMs) have demonstrated strong performance in various robot control tasks. However, their deployment in real-world applications remains constrained. Even state-of-the-art LLMs, such as GPT-o4mini, frequently produce invalid action plans that violate physical constraints, such as directing a robot to an unreachable location or causing collisions between robots. This issue primarily arises from a lack of awareness of these physical constraints during the reasoning process. To address this issue, we propose a novel framework that integrates reinforcement learning with verifiable rewards (RLVR) to incentivize knowledge of physical constraints into LLMs to induce constraints-aware reasoning during plan generation. In this approach, only valid action plans that successfully complete a control task receive positive rewards. We applied our method to two small-scale LLMs: a non-reasoning Qwen2.5-3B-Instruct and a reasoning Qwen3-4B. The experiment results demonstrate that constraint-aware small LLMs largely outperform large-scale models without constraints, grounded on both the BoxNet task and a newly developed BoxNet3D environment built using MuJoCo. This work highlights the effectiveness of grounding even small LLMs with physical constraints to enable scalable and efficient multi-robot control in complex, physically constrained environments.

 GitHub  HuggingFace

1. Introduction

Robotic control task requires controllers to find action plans given the robot's physical constraints. Conventional methods often employ planning tools, such as PDDL [1] and temporal logics [2] to find optimal plans. However, they often demand expert knowledge to convert task constraints to formal language and struggle to scale efficiently in multi-robot systems due to increased search time [3, 4, 5]. Recent advances in Large Language Models (LLMs), which excel at complex reasoning tasks like math and coding [6, 7, 8, 9], have inspired their application in robotic control. LLMs can interpret natural-language task instructions and generate valid action plans [10, 11, 12]; for instance, ChatGPT can effectively generate high-level commands such as "Robot A, move the square object to panel 2" [11, 13]. Paired with low-level execution functions that translate these commands into control signals for robots, they have proven successful in various multi-robot tasks [11, 13, 14].

However, these successes have mainly been in synthetic or constrained environments, where physical interactions are overly simplified. For example, most tasks in RocoBench have predefined all the possible valid robot interactions with the objects, largely restricting the action space for LLMs. This has led to significant issues in real-world scenarios, where LLM planners tend to violate many basic physical constraints. In particular, two important constraints are often overlooked.

*Correspondance: <jiabaoji@ucsb.edu>, <chang87@ucsb.edu>

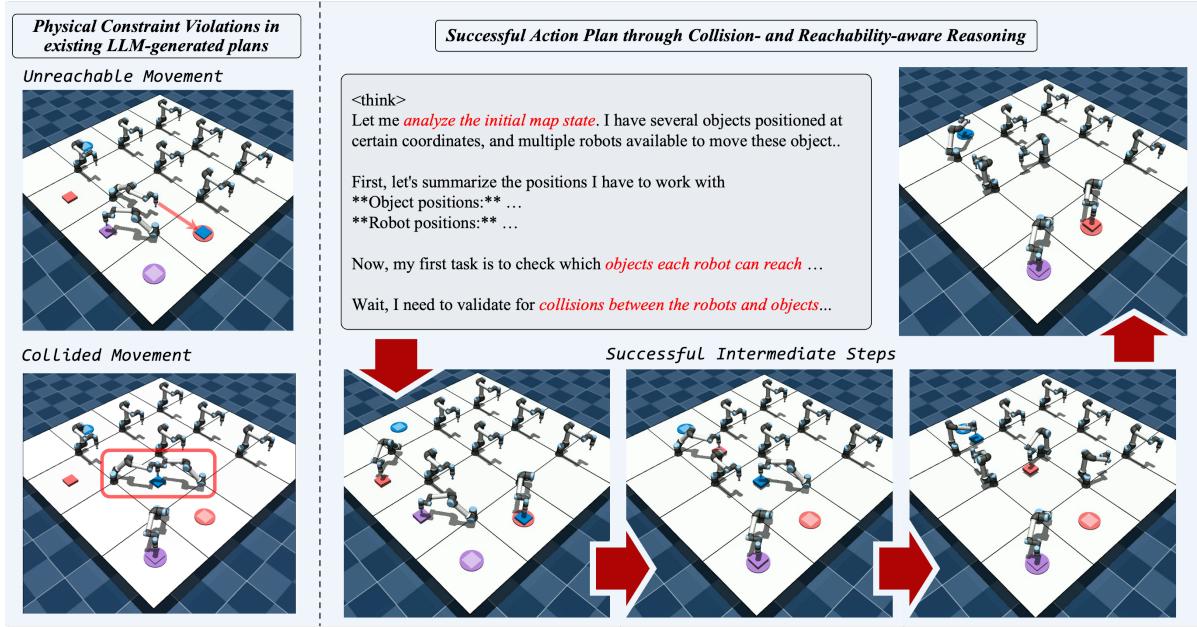


Figure 1 | Illustration of LLM-based multi-robot control. (Left) Without grounding constraint knowledge, the LLM generates action plans that result in unreachable positions or collisions. (Right) Our planner generates valid movement actions through constraint-aware reasoning (highlighted in red) that successfully completes the BoxNet task after grounding robotic constraints knowledge.

- **Reachability constraint:** LLM would direct a robot arm to an unreachable position [11, 15].
- **Collision constraint:** LLM would schedule robots to the same space, leading to collisions [13, 16].

As an example, Figure 1 (left) shows invalid actions generated by a SOTA reasoning LLM GPT-o4mini, which easily violate these constraints, leading to significant safety and feasibility concerns.

These issues highlight the imperative to equip LLM planners with the ability to understand, analyze, and adhere to basic physical constraints. However, incorporating these constraints would require strong geometric reasoning and self-reflection capabilities, particularly when the number of robots is large, which may pose nontrivial challenges to LLMs. This raises a key research question: *Can LLMs, given their current reasoning capabilities, be trained to integrate physical constraints into the planning process? If so, to what extent can they succeed?*

To study these research questions, this paper presents a novel framework to incentivize these physical constraints into LLM planners, enabling them to reason about action validity during plan generation. Specifically, we leverage reinforcement learning with verifiable rewards (RLVR) that incorporates checks for reachability, kinematic feasibility, and collision avoidance. By using binary success/failure signals derived from the robot control environment, we ensure that the LLM only receives rewards for generating physically valid plans. This fine-tuning process enables the LLM to reason about the validity during plan generation, leading to more reliable and collision-free action plans.

Our experiments on two LLMs, a non-reasoning Qwen2.5-3B-Instruct and a reasoning Qwen3-4B, have shown several encouraging findings. First, by incorporating the physical constraints into the reward, LLM planners can quickly acquire the ability to adhere to the physical constraints, thus drastically increasing the planning success rate, outperforming SOTA large-scale LLMs. For example, our best planner can achieve 0.87 and 0.53 pass@1 on two BoxNet-task multi-robot datasets, while the best baseline planner can only achieve 0.37 and 0.33 pass@1, respectively. Figure 1 (right) shows the thinking process and the generated plan by our fine-tuned LLM, which successfully solves the task without violating physical constraints. Second, our reasoning probing experiments have revealed that LLMs indeed learn to correctly identify whether the geometric constraints are satisfied or not. Finally, such capabilities acquired from RL can generalize to unseen geometric configurations, which further verifies that LLMs

learn the generic geometric reasoning skills rather than overfitting to specific geometric configurations.

In summary, the contributions of this work are as follows:

- We propose a novel framework that grounds LLMs with knowledge of action validity and collision constraints, ensuring LLM-planner-generated plans avoid unreachable positions, object collisions, or robot collisions.
- We introduce two new environments based on BoxNet task, which incorporate realistic physical constraints and serve as testbeds for evaluating LLM-based multi-robot control.
- We implement our approach on two small-scale LLMs, demonstrating that even small models like Qwen2.5-3B-Instruct and Qwen3-4B—when grounded with physical constraints—can outperform larger, state-of-the-art LLMs in complex multi-robot control tasks.

2. Method

2.1. Overview

In this section, we introduce our framework for grounding LLMs with reachability and collision awareness. Denote $\mathcal{M}_{\theta_0}(\cdot)$ as the initial LLM for performing the robot control tasks, which is capable of generating an action plan $s \sim \mathcal{M}_{\theta_0}(q; C)$ for solving the given control task described by q under a set of physical constraints C , such as the reachability of a robot arm and collision avoidance. Our goal is to fine-tune the LLM such that the generated solution s successfully moves objects to their target positions while not violating the constraints C . In the following, we first introduce the RLVR framework for grounding physical constraints in Section 2.2, then introduce our initial LLM policy warmup strategy in Section 2.3, and two different planner modes we consider in Section 2.4.

2.2. Grounding LLM with Physical Constraints through RLVR

We adopt a similar RL framework to the DeepSeek-R1 LLM [7, 17], which employs the group relative policy optimization (GRPO) algorithm [8]. Specifically, at each training step i , $i \geq 1$, we sample a group of different plans and its corresponding reasoning $\{s_1, s_2, \dots, s_G\}$ from the old LLM policy $\mathcal{M}_{\theta_{i-1}}$ for each query robot-control task q , where G is the group size. Each plan s_j is simulated in a manually implemented environment with physical constraints. The corresponding reward function $r(\cdot)$ later estimates whether it successfully completes the given task with the simulation environment feedback. Then the LLM is optimized by maximizing the following objective [8].

$$\begin{aligned} \mathcal{J}_{GRPO}(\mathcal{M}_{\theta_i}) &= \mathbb{E} \left[\left(q \sim \mathcal{D}, \{s_j\}_{j=1}^G \sim \mathcal{M}_{\theta_{i-1}}(O|q; C) \right) \right] \\ &= \frac{1}{G} \sum_{j=1}^G \left(\min \left(\frac{\mathcal{M}_{\theta}(s_j|q; C)}{\mathcal{M}_{\theta_{i-1}}(s_j|q; C)}, A_j, \text{clip} \left(\frac{\mathcal{M}_{\theta}(s_j|q; C)}{\mathcal{M}_{\theta_{i-1}}(s_j|q; C)}, 1 - \epsilon, 1 + \epsilon \right) A_j \right) \right. \\ &\quad \left. - \beta \mathbb{D}_{KL} (\mathcal{M}_{\theta_i} \| \mathcal{M}_{\theta_0}) \right), \end{aligned}$$

where \mathcal{D} denotes the training data and A_j represents the advantage, computed as the reward of each plan subtracted by the average reward within the group. Detailed definitions are in Appendix B.

Our reward function, denoted as $r(\cdot)$ largely follows the design in DeepSeek-R1 [7], with an additional plan efficiency term. Specifically,

$$r(s; q, s^*, C) = r_{\text{format}}(s) + r_{\text{execute}}(s; C) - r_{\text{efficiency}}(s; s^*),$$

$r_{\text{format}}(s) = 0.1$ if the generated solution adheres to the required thinking-then-response format and 0 otherwise. $r_{\text{execute}}(s; C) = 1$ if the simulator verifies that the plan ① accomplishes the task AND ② no physical constraints are violated, and 0 otherwise. Incorporating physical constraint checking in r_{execute} is the key mechanism to improve constraint awareness of the LLM planner. Finally,

$$r_{\text{efficiency}}(s; s^*) = \max(0, 0.1 \times (\text{len}(s) - \text{len}(s^*)))$$

penalizes the excessive length compared to a golden plan, s^* , which is obtained by an A* search algorithm (see Appendix B for search details). Additionally, we implement a minimum cap of $2 \times r_{\text{format}}(s)$ when s is a valid plan that successfully solves the task to ensure that correct plans always receive a higher reward than incorrect ones.

Table 1 | Example synthesized reasoning trace for FULLPLAN planner and REPLAN planner. We highlight three reasoning behaviors: 1. environment analysis in blue, 2. validity verification in green, and 3. efficiency consideration in orange. The environment feedback is marked in gray.

FULLPLAN Planner Reasoning:	REPLAN Planner Reasoning:
Object positions: Object 0: [0.75, 1.75] ... <think> Okay, let me analyze the given environment before coming up with a movement plan ### Evaluating Reachability: - **Robot 0** cannot reach either object ... ### Collision Checks: - No objects are placed at the same position... Let me finalize this, since there are no collisions, this plan efficiently ensures all objects are moved ...	<observation> Object positions: Object 0: [0.75, 2.75] ... </observation> <think> Alright, let's break down the situation ... <observation> Object positions: Object 1: [2.25, 2.75] ... </observation> <think> Alright, let's analyze this ...

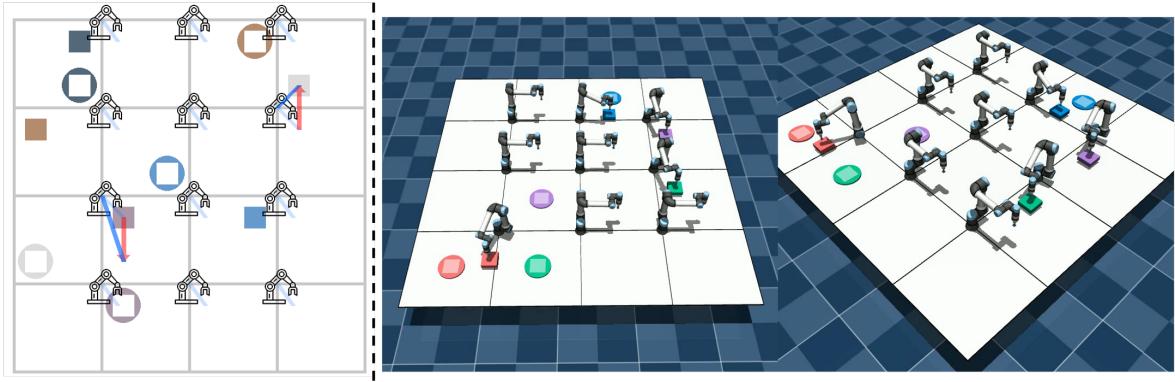


Figure 2 | (Left) An example BoxNet2D environment. The blue lines mark the robot arm, and the red lines mark the movement trajectory. (Right) An example BoxNet3D environment. Both environments require robots to collaborate to move boxes to the circle with the corresponding color.

2.3. Initial Supervised Fine-Tuning (SFT) Warmup

Prior works have shown that LLMs' initial performance on a reasoning task is crucial to RLVR training [6, 8, 9]. Since off-the-shelf LLMs often struggle with robotic control tasks, we introduce an SFT warmup to equip them with basic robot control knowledge before RL training.

The SFT data need to contain two components: ❶ a correct plan to solve a given task, and ❷ a reasoning chain that reflects a multi-step decision-making process leading to the correct plan. To synthesize such data, for each task, we first use the A* search algorithm to search for the optimal plan. Then, we pass the plan to an LLM, which is prompted to generate a reasoning process for the plan, consisting of the following three patterns:

- *Analysis of the given environment*, where the LLM assesses the current positions of robot arms and objects, e.g., *let me analyze the current situation*;
- *Validity verification*, where the LLM reasons about an arm's reachable area based on its base position and potential collision between different arms, e.g., *If Robot 0 moves, ..., it will collide*;
- *Efficiency considerations*, where the LLM evaluates whether multiple movements can be parallelized to improve the plan efficiency.

Table 1 shows the example reasoning chain synthesized by GPT4o-mini, where the three patterns are rendered in different colors. Appendix D shows the full prompt for our reasoning synthesis.

2.4. Two Planners: FULLPLAN Planner and REPLAN Planner

We consider two different LLM-based planners in this work. The first planner, referred to as **FULLPLAN**, involves the LLM directly generating the entire plan that may take multiple execution steps for solving a task based on the initial positions of all objects and robot poses in the environment. The second planner, denoted as **REPLAN**, generates one step at a time, observing the updated object positions from the environment (appended to its context) before generating the next step. This allows the planner to evolve dynamically as the environment changes through multiple execution steps. Table 1 provides the example planning processes for two different LLM-based planners for the same initial environment. We highlight that the **REPLAN** planner has access to multiple intermediate observations of the environment, while the **FULLPLAN** planner only sees the initial environment.

3. BoxNet-Based Multi-Robot Environments

In this work, we primarily experiment with BoxNet task [11], where multiple robots collaborate to move objects across different cells to targeted locations in a fixed grid map. This section details two environments we developed, a modified BoxNet2D environment and the newly developed BoxNet3D environment, both equipped with realistic physical constraint checks.

BoxNet2D. Figure 2 (left) shows a BoxNet2D environment. In this environment, robot arms are placed at a corner of a grid environment. Each arm can reach its neighboring grids for picking and placing objects. Unlike the previous BoxNet environment that predefined all valid robot arm actions, we allow LLMs to generate spatial coordinates directly, significantly expanding the action space. For example, the action “*Robot0 Move (1.25, 1.25) → (1.75, 1.75), False*” moves *Robot0*’s arm to (1.75, 1.75) without picking up an object. In contrast, “*Robot1 Move (2.25, 1.75) → (1.25, 1.25), True*” indicates *Robot1* picking up the object at (1.75, 1.25) and moving it to (2.25, 1.75).

We pre-define four points within each grid, *e.g.*, (0.25, 0.25), (0.25, 0.75), (0.75, 0.25), (0.75, 0.75), for object placement and robot arm moving, and later we will show that the fine-tuned LLM can generalize to other points in experiments. Three physical constraints are implemented: ① *reachability verification*, which checks whether the target position of a robot is unreachable. ② *robot collision detection*, which checks whether the movement trajectory of different arms intersects with each other, or one robot’s movement trajectory intersects with another robot arm, leading to a collision. ③ *object collision detection*, which checks whether two objects are placed at the same spatial coordinates during the plan execution. Example invalid actions of BoxNet2D are provided in Appendix C.

BoxNet3D. Figure 2 (right) shows a BoxNet3D environment. In this environment, we employ the UR5e robot arm as the basic robot arm¹. Similar to the 2D environment, the goal is to move colored boxes into corresponding circles of the same color with the fewest actions. Each robot arm’s base is fixed and moves its arm around to reach different grids for object picking and placement. We employ an RRT planner implemented by RoCoBench [13] for low-level control signal generation given LLM-generated coordinates for arm position movement². We employ Mujoco as the simulation engine [18], which provides the arm reachability check and collision detection mechanism.

4. Experiment

In this section, we conduct empirical experiments on the two BoxNet-based environments to assess the effectiveness of our method. We first present the experiment setup in Section 4.1 and then the experiment results in Section 4.2, followed by additional ablation studies in Section 4.3.

¹<https://www.universal-robots.com/products/ur5e/>

²The RRT planner is adapted from the implementation in RoCoBench official code base (<https://github.com/MandiZhao/robot-collab/blob/main/rocobench/rrt.py>)

Table 2 | Performance of different LLM planners on BoxNet2D and BoxNet3D . For each model, we report the results for FULLPLAN planner and REPLAN planner side-by-side (FULLPLAN / REPLAN).

Model	BoxNet2D			BoxNet3D		
	Success ↑	StepDiff. ↓	Para. ↑	Success ↑	StepDiff. ↓	Para. ↑
Search Algorithm						
A*	1	0	2.24	1	0	2.14
LLMs without constraint knowledge grounding						
GPT-4omini	0.06 / 0.05	2.35 / 0.14	1.17 / 1.15	0.07 / 0.06	0.79 / 0.45	1.03 / 1.08
GPT-4o	0.12 / 0.11	2.14 / 0.13	1.15 / 1.22	0.10 / 0.12	0.23 / 0.68	1.35 / 1.11
GPT-o4mini	0.37 / 0.35	0.24 / -0.31	1.58 / 1.87	0.11 / 0.33	0.14 / 1.21	1.45 / 1.53
Qwen2.5-3B-Inst	0.0 / 0.0	— / —	— / —	0.08 / 0.0	0.20 / —	1.40 / —
Qwen2.5-7B-Inst	0.02 / 0.02	1.45 / 0.31	1.20 / 1.23	0.05 / 0.08	0.41 / 0.35	1.13 / 1.07
QwQ-32B	0.04 / 0.07	0.35 / 0.17	1.12 / 1.21	0.07 / 0.15	0.24 / -0.09	1.08 / 1.31
Qwen3-4B	0.14 / 0.13	0.23 / 0.14	1.29 / 1.29	0.15 / 0.11	0.07 / 0.31	1.17 / 1.14
Qwen3-8B	0.18 / 0.15	-0.23 / -0.34	1.24 / 1.31	0.17 / 0.13	-0.02 / 0.09	1.22 / 1.21
Qwen3-14B	0.19 / 0.21	-0.31 / -0.24	1.34 / 1.41	0.10 / 0.14	0.17 / 1.37	1.34 / 1.35
Qwen3-32B	0.11 / 0.14	0.17 / -0.03	1.24 / 1.12	0.14 / 0.17	0.09 / 0.04	1.27 / 1.09
LLMs with grounded constraint knowledge						
Qwen2.5-3B-SFT	0.34 / 0.30	0.11 / -0.04	1.51 / 1.39	0.27 / 0.39	-0.07 / -0.05	1.27 / 1.39
Qwen2.5-3B-RL	0.58 / 0.68	-0.65 / 0.23	1.53 / 1.50	0.42 / 0.48	-0.15 / -0.14	1.33 / 1.49
Qwen3-4B-SFT	0.45 / 0.31	-0.12 / -0.15	1.92 / 1.35	0.37 / 0.43	0.09 / -0.11	1.32 / 1.48
Qwen3-4B-RL	0.87 / 0.75	-0.84 / -0.64	1.73 / 1.64	0.45 / 0.53	-0.25 / -0.29	1.39 / 1.56

4.1. Experiment Setup

Dataset generation. We create environments with various map sizes and object initial and target positions for both BoxNet2D and BoxNet3D . Specifically, for BoxNet2D , we use map sizes ranging from 2×2 to 6×6 and 1 to 5 objects, resulting in 55,000 training and 250 testing environments. For BoxNet3D , we use map sizes from 2×2 to 4×4 with 1 to 3 objects, yielding 1,800 training and 160 testing environments. The object position is randomly sampled from the pre-defined points, while the robot arms are evenly placed at the grid joints to ensure that all grids in the map can be reached. For each randomly sampled environment, the manually implemented A* search algorithm verifies that a valid solution action plan exists. Detailed dataset statistics are summarized in Appendix B.2.

Evaluation metric. We evaluate LLM-based planners mainly from two perspectives: ① *Success*, the proportion of generated plans that solve given robotic tasks, measured by *pass@1* over four trials per environment; and ② *StepDiff.*, the difference in number of steps between successful plans and the best plan among A* solutions. We also report *Para.*, the maximum number of robots operating in parallel in any intermediate step of a successful plan.

Baseline LLMs. We mainly compare with off-the-shelf LLMs. To ensure comprehensive coverage of existing LLMs, our evaluation includes both reasoning and non-reasoning models, closed-source and open-source ones across different parameter scales. Specifically, we consider closed-source LLMs, GPT-4o, GPT-4o-mini, and GPT-o4-mini. On the open-source side, we include Qwen-2.5 and Qwen3 series, with parameter sizes ranging from 3B to 32B.

Training details. We use two base LLMs: a non-reasoning LLM Qwen-2.5-3B-Instruct and a reasoning LLM Qwen3-4B. For SFT warm-up, we use a learning rate of $1e-5$ for Qwen-2.5-3B-Instruct and $3e-5$ for Qwen-3-4B with the AdamW optimizer [19]. Training runs for 10 epochs on FULLPLAN and 5 epochs on REPLAN . RLVR training uses a fixed $1e-6$ learning rate for 200 steps with the GRPO algorithm [20]. Batch sizes are 256 (group size 8) for BoxNet2D and 64 for BoxNet3D . Following prior work [6, 9], we set $\beta = 0$ in the GRPO loss. We use the VeRL framework [21], and run all experiments on 2×8 NVIDIA H100 GPUs.

4.2. Experimental Results

Grounding empowers small-scale LLMs to outperform larger ones. We first evaluate the grounded LLM planner performance in Table 2. The LLMs with physical constraints knowledge grounded through SFT warmup and further RL training are denoted by the suffix *-SFT* and *-RL*, respectively. We highlight three observations: First, grounding constraint knowledge significantly boosts planning performance, enabling 3B and 4B LLMs to achieve higher success against larger ones. For example, Qwen3-4B-RL FULLPLAN planner achieves 0.87 success rate, 0.5 higher than the best baseline GPT-o4mini. Second, grounded LLM planners produce more efficient plans than the A* search algorithm on solved tasks. For example, Qwen3-4B-RL has 0.84 fewer steps than the ground-truth plan from our A* implementation, showing a strong reasoning ability and also echoes the findings in prior works that compare LLM planners with A* on Sudoku [22, 23]. Third, planner performance on BoxNet3D is generally worse than on BoxNet2D. This suggests that, although we applied multiple feasibility checks in BoxNet2D, some physical constraints remain missing. The BoxNet3D environment uses a more advanced simulation engine and thus exposes more limitations. These results underscore the importance of developing realistic robotic environments that capture real-world complexity for future LLM-based robotic control research.

Figure 3 visualizes planner performance against numbers of boxes for BoxNet2D. We highlight that RL-trained planners better preserve performance when task complexity increases. For example, the performance gap between Qwen3-4B-SFT and Qwen3-4B-RL grows from 0.13 to 0.53 when the number of boxes increases from 1 to 5 for BoxNet2D, highlighting better scalability of RL planners.

RL planners generalize better to unseen environments. To measure how the planners’ reasoning ability generalizes, we evaluate the planners’ performance in two unseen variants of BoxNet2D test data: ① Random robot layout, denoted as RANDROB, where the robot positions are randomly assigned on the grid joints in maps ranging from 2×2 to 5×5 . We ensure that all testing data are solvable, which means every box can reach its target position via robot movement. ② Unseen coordinates, *i.e.*, NEWCOORD, where the initial and target position coordinates of all objects in BoxNet2D test set are perturbed by a random offset ($\Delta x, \Delta y$) $\sim \mathcal{U}([-0.2, 0.2]^2)$. Example data are shown in Appendix C.

Table 3 | Planning performance generalization on unseen BoxNet2D environments.

Model	RANDROB		NEWCOORD	
	Success \uparrow	StepDiff. \downarrow	Success \uparrow	StepDiff. \downarrow
FULLPLAN Planner				
Qwen2.5-3B-SFT	0.39	0.12	0.32	0.21
Qwen2.5-3B-RL	0.58	-0.04	0.55	-0.32
Qwen3-4B-SFT	0.48	0.03	0.43	-0.03
Qwen3-4B-RL	0.79	-0.40	0.87	-0.39
REPLAN Planner				
Qwen2.5-3B-SFT	0.39	0.23	0.33	0.09
Qwen2.5-3B-RL	0.71	1.24	0.68	1.04
Qwen3-4B-SFT	0.41	-0.03	0.37	0.15
Qwen3-4B-RL	0.75	-0.15	0.69	0.09

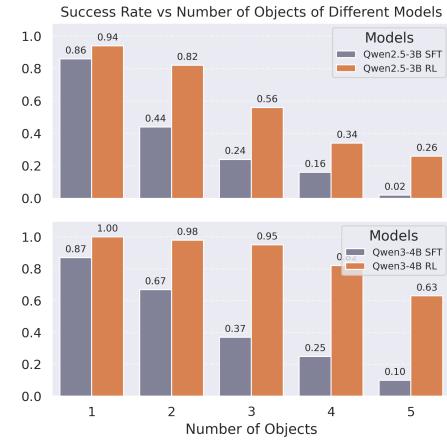


Figure 3 | Success rate against number of boxes in the BoxNet2D environment.

Table 3 reports the performance of our grounded planners on two unseen environments. We highlight that the RL-trained planners consistently outperform SFT ones while maintaining plan efficiency. For example, Qwen3-4B-RL FULLPLAN planner achieves 0.87 success rate on NEWCOORD, 0.44 better than the SFT variant, showing better generalization of the reasoning capability. This observation also aligns with previous RL for LLM works [12, 24].

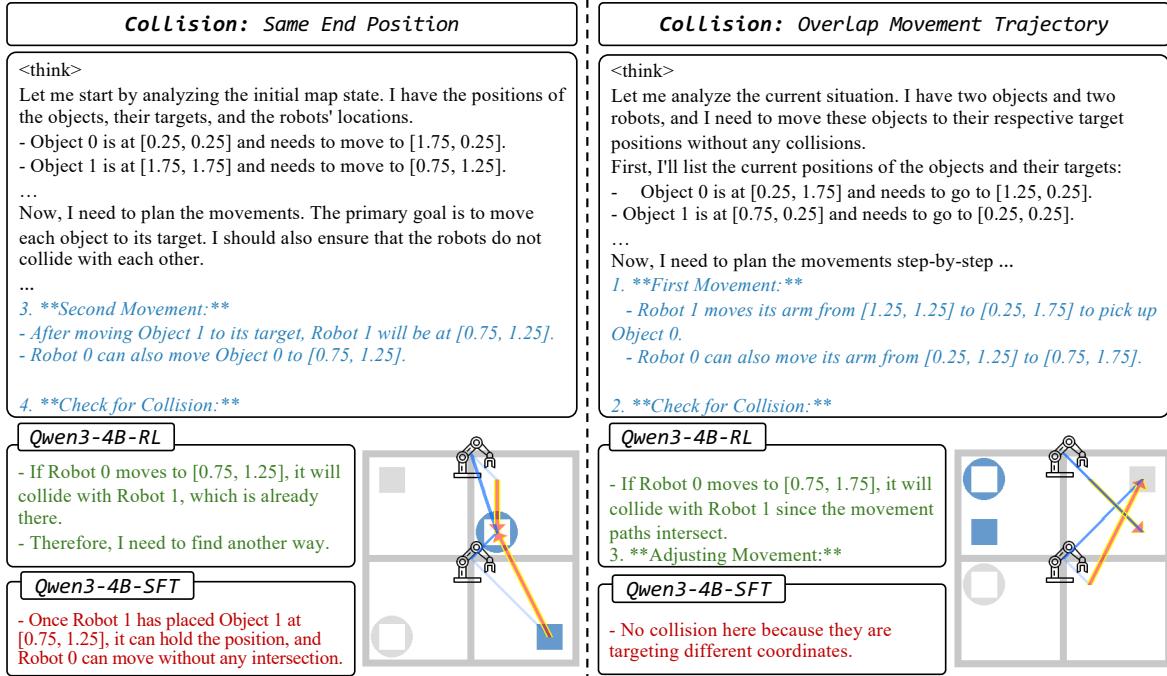


Figure 4 | Example reasoning trace generated by grounded FULLPLAN planners. Invalid action plans are manually inserted into the trace history and highlighted in blue. The correct continuations that identify and correct these errors are shown in green, while incorrect continuations are shown in red. The bottom figure visualizes the collision between two movements. RL planner better detects errors.

Reasoning behavior change after RL. Previous results have shown that RL training significantly improves the planner’s planning ability. In this section, we analyze in more detail how the reasoning behavior of the LLM-based planners changes before and after RL fine-tuning.

Given the critical role of reachability checks and collision checks in generating successful action plans, we prompt GPT-4o to count the number of these checks in the reasoning traces produced by our FULLPLAN planners across three BoxNet2D environment variants. The prompts used are provided in Appendix D. As shown in Table 4, the RL planners perform more reachability checks (Rea.) and collision checks (Col.) than the initial SFT planner. These checks help ensure the feasibility of action plans and lead to a large improvement in success rate. This observation suggests that RL training helps the LLM better understand the importance of these checks and use them more consistently.

We also conduct a qualitative analysis to verify their reasoning ability by injecting error steps into the trace. Specifically, we insert an invalid action that would lead to collision into the intermediate reasoning steps. To test whether the planner can recognize and correct such errors, we append the phrase “Collision Check” to the perturbed trace to trigger verification. Figure 4 shows the full examples, where the injected invalid actions are highlighted in blue. The LLM’s continuation is marked in green if it identifies and corrects the error, and in red if it fails. We find that the RL-trained planner successfully finds the error and traces the issue to same target position and overlapped movement paths. This suggests that RL helps build better physical constraints-aware reasoning.

4.3. Ablation Study

In this section, we explore the design choices for our framework on the BoxNet2D environment, focusing on: ① How does SFT warmup affect final planner performance? ② Is the textual thinking necessary for planner performance? ③ How does the efficiency penalty affect the planner’s behavior?

Initial LLM policy matters. Figure 5 shows how the training reward evolves over the first 40 steps for different initial FULLPLAN planners. For the original Qwen2.5-3B-Instruct and Qwen3-4B models, we observe a sharp reward increase of about 0.1 (the format reward) within the first 10 steps, which indicates that they quickly learn to produce answers in required output format. However, after this initial gain, the reward plateaus, indicating limited additional learning to improve planning capability. In contrast, the LLM with SFT warmup shows a consistent increasing reward trend. This suggests that the SFT warmup helps build a strong foundation, allowing it to continue learning and optimize effectively in RL training.

Textual reasoning improves planning performance.

To assess the role of textual reasoning for LLM planners, we perform an ablation study with our SFT and RL pipeline on Qwen-2.5-3B-Instruct. In this experiment, we train a planner without the synthesized thinking steps, *i.e.*, it generates only the final action plan with no textual thinking. As shown in Table 5, removing intermediate reasoning leads to a notable performance drop: success rates fall from 0.34 to 0.26 for SFT planner, and from 0.58 to 0.39 for RL planner. This highlights the importance of textual thinking for LLM planners.

Efficiency penalty in reward improves plan efficiency. We observed a surprising finding that RL-trained LLM planners produce more efficient plans than those generated by our hand-crafted A* search algorithm in Table 2, which is likely due to the efficiency penalty term in the reward function. To further understand its role, we conduct an ablation study on $r_{\text{efficiency}}$. Starting from the same initial LLM policy Qwen-2.5-3B-SFT, we perform RL training without the efficiency penalty for FULLPLAN planner. Table 5 presents the results. While both RL-finetuned LLM largely improve the success rate, their plan efficiency differs significantly. The planner trained with $r_{\text{efficiency}}$ produces plans that are 2.09 steps shorter. In contrast, the parallelism drops close to 1 without the penalty in RL. These results underscore the importance of efficiency penalty in reward.

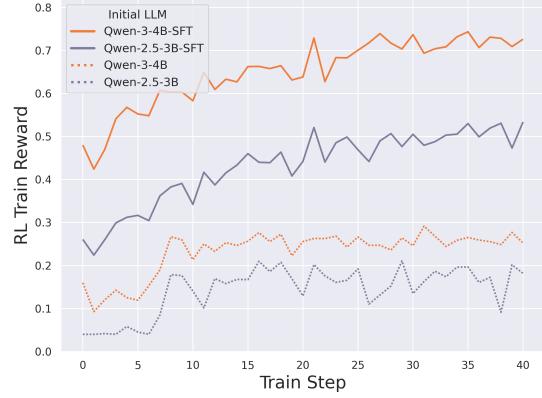


Figure 5 | Training reward trajectory in first 40 steps for different initial LLMs.

Table 5 | Impact of ablating thinking and $r_{\text{efficiency}}$ on BoxNet2D performance.

Model	BoxNet2D		
	Success ↑	StepDiff. ↓	Para. ↑
Qwen2.5-SFT	0.34	0.11	1.51
– thinking	0.26	0.05	1.35
Qwen2.5-RL	0.58	-0.65	1.53
– thinking	0.39	-0.07	1.43
– $r_{\text{efficiency}}$	0.52	1.44	1.07

5. Related Work

Robotic planning and control with LLMs. Robotic planning and control is a complex task that requires high-level planning under various physical constraints. Traditional approaches typically translate task goals and physical constraints into formal logic specifications, such as Temporal Logic or PDDL [1, 2], and solve them using constraint solvers. More recently, LLMs have been applied to robotic control due to their strong reasoning capabilities and better scalability compared to constraint solvers. For example, some works use LLMs to choose actions from predefined motion primitives [25, 26, 27], while others treat code as an intermediate representation for control [4, 5, 10, 28, 29, 30]. Hybrid approaches, such as

AutoTAMP [3] and Text2Motion [31], combine traditional planning tools with LLMs for action planning. Another series of works employs multi-LLM discussion for robotic tasks [11, 13, 32, 33, 34]. While these methods show promising results, many of them overly simplify physical constraints, limiting their real-world applicability. In this work, we demonstrate that even SOTA LLMs struggle under realistic physical constraints, and further introduce a novel approach that grounds smaller LLMs with this constraint knowledge, which largely improves performance.

Reinforcement learning with verifiable rewards for LLM reasoning. Reinforcement learning (RL) has demonstrated significant promise in enhancing the reasoning capabilities of large language models (LLMs) across a wide range of domains, including mathematics [17, 20, 35], code generation [9, 6, 36], and complex multi-agent systems [37, 38, 39, 40]. A common paradigm involves training LLMs to optimize for a verifiable reward, such as the correctness of a math solution or whether the generated code passes unit tests, using RL training. Many previous works show that the RL training process vastly improves LLM reasoning [12, 39, 41, 42, 43]. The improvement is often accompanied by emergent reasoning behaviors such as feasibility checks and self-reflection, which are difficult to elicit through supervised fine-tuning alone [12, 44, 45]. In this work, we extend the RLVR to robotic control, with a focus on grounding LLMs with knowledge of physical constraints. Our method leverages RLVR to teach LLMs to reason under the constraints inherent in robotic planning tasks, such as collision avoidance, reachability accordance, and goal satisfaction. By integrating physical verification signals into the training process, the model learns to internalize these constraints as part of its reasoning process. This grounding leads to more robust and reliable control planning decisions in downstream robotic applications.

6. Conclusion

In this paper, we present a novel framework that grounds LLMs with physical constraint knowledge, such as robot arm reachability and collision avoidance. By incorporating these constraints, LLMs are able to reason more effectively about action feasibility and generate efficient and physically viable action plans. To evaluate our approach, we developed two BoxNet-based multi-robot environments, BoxNet2D and BoxNet3D, both equipped with action feasibility checks. Experiments show that even small-scale LLMs at 3B and 4B parameter size, when grounded with constraint knowledge, significantly outperform larger SOTA LLMs. Additional experiments on reasoning behavior and generalization further confirm that our models learn constraint-aware reasoning rather than simply overfitting to training data. We discuss limitations and potential societal impacts in Appendix A.

References

- [1] Maria Fox and Derek Long. Pddl2. 1: An extension to pddl for expressing temporal planning domains. *Journal of artificial intelligence research*, 20:61–124, 2003.
- [2] E Allen Emerson. Temporal and modal logic. In *Formal models and semantics*, pages 995–1072. Elsevier, 1990.
- [3] Yongchao Chen, Jacob Arkin, Charles Dawson, Yang Zhang, Nicholas Roy, and Chuchu Fan. Autotamp: Autoregressive task and motion planning with llms as translators and checkers. In *2024 IEEE International conference on robotics and automation (ICRA)*, pages 6695–6702. IEEE, 2024.
- [4] Yongchao Chen, Yilun Hao, Yang Zhang, and Chuchu Fan. Code-as-symbolic-planner: Foundation model-based robot planning via symbolic code generation. *arXiv preprint arXiv: 2503.01700*, 2025.
- [5] Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, et al. Inner monologue: Embodied reasoning through planning with language models. *arXiv preprint arXiv:2207.05608*, 2022.
- [6] Michael Luo, Sijun Tan, Roy Huang, Ameen Patel, Alpay Ariyak, Qingyang Wu, Xiaoxiang Shi, Rachel Xin, Colin Cai, Maurice Weber, Ce Zhang, Li Erran Li, Raluca Ada Popa, and Ion Stoica. Deepcoder: A fully open-source 14b coder at o3-mini level.

- [7] DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv: 2501.12948*, 2025.
- [8] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv: 2402.03300*, 2024.
- [9] Jiawei Liu and Lingming Zhang. Code-r1: Reproducing r1 for code with reliable rewards. 2025.
- [10] Yue Meng, Fei Chen, Yongchao Chen, and Chuchu Fan. Audere: Automated strategy decision and realization in robot planning and control via llms. *arXiv preprint arXiv: 2504.03015*, 2025.
- [11] Yongchao Chen, Jacob Arkin, Yang Zhang, Nicholas Roy, and Chuchu Fan. Scalable multi-robot collaboration with large language models: Centralized or decentralized systems? In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4311–4317. IEEE, 2024.
- [12] Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V. Le, Sergey Levine, and Yi Ma. Sft memorizes, rl generalizes: A comparative study of foundation model post-training. *arXiv preprint arXiv: 2501.17161*, 2025.
- [13] Zhao Mandi, Shreeya Jain, and Shuran Song. Roco: Dialectic multi-robot collaboration with large language models, 2023.
- [14] Dawei Sun, Jingkai Chen, Sayan Mitra, and Chuchu Fan. Multi-agent motion planning from signal temporal logic specifications. *IEEE Robotics and Automation Letters*, 7(2):3451–3458, 2022.
- [15] Songyuan Zhang, Oswin So, Kunal Garg, and Chuchu Fan. Gcbf+: A neural graph control barrier function framework for distributed safe multi-agent control. *IEEE Transactions on Robotics*, 2025.
- [16] Joshua Jones, Oier Mees, Carmelo Sferrazza, Kyle Stachowicz, Pieter Abbeel, and Sergey Levine. Beyond sight: Finetuning generalist robot policies with heterogeneous sensors via language grounding. *arXiv preprint arXiv: 2501.04693*, 2025.
- [17] Z. Z. Ren, Zhihong Shao, Junxiao Song, Huajian Xin, Haocheng Wang, Wanjia Zhao, Liyue Zhang, Zhe Fu, Qihao Zhu, Dejian Yang, Z. F. Wu, Zhibin Gou, Shirong Ma, Hongxuan Tang, Yuxuan Liu, Wenjun Gao, Daya Guo, and Chong Ruan. Deepseek-prover-v2: Advancing formal mathematical reasoning via reinforcement learning for subgoal decomposition. *arXiv preprint arXiv: 2504.21801*, 2025.
- [18] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. pages 5026–5033, 2012.
- [19] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. *International Conference on Learning Representations*, 2017.
- [20] Daya Guo, Dejian Yang, Haowei Zhang, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv: 2501.12948*, 2025.
- [21] Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint arXiv: 2409.19256*, 2024.
- [22] Lucas Lehnert, Sainbayar Sukhbaatar, Paul Mcvay, Michael Rabbat, and Yuandong Tian. Beyond a*: Better planning with transformers via search dynamics bootstrapping. *arXiv preprint arXiv: 2402.14083*, 2024.
- [23] DiJia Su, Sainbayar Sukhbaatar, Michael Rabbat, Yuandong Tian, and Qinqing Zheng. Dualformer: Controllable fast and slow thinking by learning with randomized reasoning traces. In *The Thirteenth International Conference on Learning Representations*, 2024.
- [24] Noam Razin, Zixuan Wang, Hubert Strauss, Stanley Wei, Jason D. Lee, and Sanjeev Arora. What makes a reward model a good teacher? an optimization perspective. *arXiv preprint arXiv: 2503.15477*, 2025.

- [25] Lin Guan, Karthik Valmeekam, Sarath Sreedharan, and Subbarao Kambhampati. Leveraging pre-trained large language models to construct and utilize world models for model-based task planning. *Advances in Neural Information Processing Systems*, 36:79081–79094, 2023.
- [26] Marta Skreta, Naruki Yoshikawa, Sebastian Arellano-Rubach, Zhi Ji, Lasse Bjørn Kristensen, Kourosh Darvish, Alán Aspuru-Guzik, Florian Shkurti, and Animesh Garg. Errors are useful prompts: Instruction guided task programming with verifier-assisted iterative prompting. *arXiv preprint arXiv:2303.14100*, 2023.
- [27] João Loula, Benjamin LeBrun, Li Du, Ben Lipkin, Clemente Pasti, Gabriel Grand, Tianyu Liu, Yahya Emara, Marjorie Freedman, Jason Eisner, Ryan Cotterell, Vikash Mansinghka, Alexander K. Lew, Tim Vieira, and Timothy J. O’Donnell. Syntactic and semantic control of large language models via sequential monte carlo. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net, 2025.
- [28] Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. Code as policies: Language model programs for embodied control. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9493–9500. IEEE, 2023.
- [29] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, K. Gopalakrishnan, Karol Hausman, Alexander Herzog, Daniel Ho, Jasmine Hsu, Julian Ibarz, Brian Ichter, A. Irpan, Eric Jang, Rosario M Jauregui Ruano, Kyle Jeffrey, Sally Jesmonth, N. Joshi, Ryan C. Julian, Dmitry Kalashnikov, Yuheng Kuang, Kuang-Huei Lee, S. Levine, Yao Lu, Linda Luu, Carolina Parada, P. Pastor, Jornell Quiambao, Kanishka Rao, Jarek Rettinghouse, D. Reyes, P. Sermanet, Nicolas Sievers, Clayton Tan, Alexander Toshev, Vincent Vanhoucke, F. Xia, Ted Xiao, Peng Xu, Sichun Xu, and Mengyuan Yan. Do as i can, not as i say: Grounding language in robotic affordances. *Conference on Robot Learning*, 2022.
- [30] Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. Progprompt: Generating situated robot task plans using large language models. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11523–11530. IEEE, 2023.
- [31] Kevin Lin, Christopher Agia, Toki Migimatsu, Marco Pavone, and Jeannette Bohg. Text2motion: From natural language instructions to feasible plans. *arXiv preprint arXiv: 2303.12153*, 2023.
- [32] Yang Zhang, Shixin Yang, Chenjia Bai, Fei Wu, Xiu Li, Zhen Wang, and Xuelong Li. Towards efficient llm grounding for embodied multi-agent collaboration. *arXiv preprint arXiv:2405.14314*, 2024.
- [33] Xudong Guo, Kaixuan Huang, Jiale Liu, Wenhui Fan, Natalia Vélez, Qingyun Wu, Huazheng Wang, Thomas L Griffiths, and Mengdi Wang. Embodied llm agents learn to cooperate in organized teams. *arXiv preprint arXiv:2403.12482*, 2024.
- [34] Zhixuan Shen, Haonan Luo, Kexun Chen, Fengmao Lv, and Tianrui Li. Enhancing multi-robot semantic navigation through multimodal chain-of-thought score collaboration. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 14664–14672, 2025.
- [35] Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He. Simplerl-zoo: Investigating and taming zero reinforcement learning for open base models in the wild. *arXiv preprint arXiv: 2503.18892*, 2025.
- [36] OpenAI. Competitive programming with large reasoning models. *arXiv preprint arXiv: 2502.06807*, 2025.
- [37] Bowen Jin, Hansi Zeng, Zhenrui Yue, Dong Wang, Hamed Zamani, and Jiawei Han. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *arXiv preprint arXiv: 2503.09516*, 2025.
- [38] Joykirat Singh, Raghav Magazine, Yash Pandya, and Akshay Nambi. Agentic reasoning and tool integration for llms via reinforcement learning. *arXiv preprint arXiv:2505.01441*, 2025.

- [39] Peiyuan Feng, Yichen He, Guanhua Huang, Yuan Lin, Hanchong Zhang, Yuchen Zhang, and Hang Li. Agile: A novel reinforcement learning framework of llm agents. *Neural Information Processing Systems*, 2024.
- [40] OpenAI. Introducing operator, 2024. Accessed: 2025-03-23.
- [41] Jiazhen Pan, Che Liu, Junde Wu, Fenglin Liu, Jiayuan Zhu, Hongwei Bran Li, Chen Chen, Cheng Ouyang, and Daniel Rueckert. Medvilm-r1: Incentivizing medical reasoning capability of vision-language models (vlms) via reinforcement learning. *arXiv preprint arXiv:2502.19634*, 2025.
- [42] Haozhan Shen, Peng Liu, Jingcheng Li, Chunxin Fang, Yibo Ma, Jiajia Liao, Qiaoli Shen, Zilun Zhang, Kangjia Zhao, Qianqian Zhang, et al. Vlm-r1: A stable and generalizable r1-style large vision-language model. *arXiv preprint arXiv:2504.07615*, 2025.
- [43] Bairu Hou, Yang Zhang, Jiabao Ji, Yujian Liu, Kaizhi Qian, Jacob Andreas, and Shiyu Chang. Thinkprune: Pruning long chain-of-thought of llms via reinforcement learning. 2025.
- [44] Eric Zelikman, Georges Raif Harik, Yijia Shao, Varuna Jayasiri, Nick Haber, and Noah Goodman. Quiet-star: Language models can teach themselves to think before speaking. In *First Conference on Language Modeling*, 2024.
- [45] Arian Hosseini, Xingdi Yuan, Nikolay Malkin, Aaron Courville, Alessandro Sordoni, and Rishabh Agarwal. V-star: Training verifiers for self-taught reasoners. *arXiv preprint arXiv:2402.06457*, 2024.
- [46] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [47] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.

A. Limitations and Societal Impacts

Our work introduces a novel framework to ground LLMs with physical constraint knowledge for robot control tasks, significantly enhancing their ability to reason about action feasibility during plan generation. This leads to substantial improvements in planning performance. We validate the effectiveness of our approach through experiments in two BoxNet-based multi-robot environments on two small-scale LLMs. However, this work has two limitations: ① Our experiments are limited to the BoxNet task due to the high implementation overhead required for other robot control environments. Extending our framework to additional physical constraint-sensitive tasks remains an important direction for future work. ② The RL training is conducted at a limited scale due to computational constraints. Unlike typical RL setups in math and coding domains that allow for training over multiple epochs [9, 6], our training is restricted to just one or two epochs. Despite this limitation, our experimental results already demonstrate strong reasoning capabilities on robotic tasks.

Our work aims to advance the integration of LLMs into robotic control planning, which has many promising societal benefits. By enabling LLMs to better understand and operate within physical constraints, LLMs can help build safer, more reliable, and more efficient multi-robot systems. This can potentially enhance robotic automation domains heavily involving many robots. In particular, improved planning performance can reduce operational errors and increase productivity. However, as with any deployment of AI in real-world decision-making systems, there are potential risks if the planners are deployed with a dangerous purpose. Future extensions of this work should also consider robustness to adversarial scenarios to ensure responsible real-world integration.

B. Additional Implementation Details

In this section, we provide more implementation details including: the RL training algorithm (Section B.1), the implementation of BoxNet2D, BoxNet3D and dataset statistics (Section B.2), and the A* search algorithm for data generation (Section B.3).

B.1. GRPO Algorithm

GRPO [8], or group relative policy optimization, is a variant of PPO algorithm [46] proposed for LLM RL. In this section, we briefly outline the GRPO algorithm and refer readers to the original paper [7] for more details.

As we mentioned in the main paper, given an initial LLM policy \mathcal{M}_{θ_0} and a train dataset \mathcal{D} , the GRPO loss is defined as follows:

$$\begin{aligned} \mathcal{J}_{GRPO}(\mathcal{M}_{\theta_i}) &= \mathbb{E} \left[\left(\mathbf{q} \sim \mathcal{D}, \{s_j\}_{j=1}^G \sim \mathcal{M}_{\theta_{i-1}}(\mathbf{o}|\mathbf{q}; C) \right) \right] \\ &= \frac{1}{G} \sum_{j=1}^G \left(\min \left(\frac{\mathcal{M}_{\theta}(s_j|\mathbf{q}; C)}{\mathcal{M}_{\theta_{i-1}}(s_j|\mathbf{q}; C)}, A_j, \text{clip} \left(\frac{\mathcal{M}_{\theta}(s_j|\mathbf{q}; C)}{\mathcal{M}_{\theta_{i-1}}(s_j|\mathbf{q}; C)}, 1 - \epsilon, 1 + \epsilon \right) A_j \right) \right. \\ &\quad \left. - \beta D_{KL} (\mathcal{M}_{\theta_i} \| \mathcal{M}_{\theta_0}) \right), \end{aligned}$$

where i denotes the train step, G denotes the group size, \mathbf{q} denotes a textual query describing a robotic control task, C denotes the constraints in text, A_j denotes the advantage for j -th rollout s_j . The definition of advantage is:

$$A_j = \frac{r_j - \text{mean}(\mathbf{r})}{\text{std}(\mathbf{r})},$$

given the reward $\mathbf{r} = \{r_1, \dots, r_G\}$ for all LLM rollouts to the query task \mathbf{q} , following the Generalized Advantage Estimation (GAE) [47].

B.2. BoxNet Environment Implementation and Statistics

We implement two different environments based on the BoxNet task, which involves multiple robots in a grid map and collaborating to move objects to the corresponding target positions. Both environments are implemented in Python.

Table 6 | Dataset statistics of BoxNet2D and BoxNet3D. The average steps to complete and parallelism are all based on the optimal plans generated by our manually implemented A* algorithm.

Dataset	Sample No.	Avg. Step	Para.
BoxNet2D-train	55000	8.13	1.73
BoxNet2D-test	250	8.32	1.75
RANDROB	200	7.06	1.49
NEWCOORD	250	8.59	1.77
BoxNet3D-train	1800	6.27	1.89
BoxNet3D-test	160	5.62	1.88

BoxNet2D For BoxNet2D, we manually implement the feasibility check by calculating the relative geometric position of robot arms and objects. The map size ranges from 2×2 to 6×6 , and the object number ranges from 1 to 5. For each unique map configuration, *i.e.*, a tuple of map width, height, and the object number, we randomly generate at most 150 different object initial and target positions to construct the unique environments in train dataset. The testing data consists of the square maps with the width ranging from 2 to 6, and the object number ranges from 1 to 5. We generate at most 10 unique environments to construct the test dataset. The detailed dataset statistics are summarized in Table 6.

In the unseen environment transfer experiment, we generate two variants of BoxNet2D test set: RANDROB, where the robot position is not evenly placed at the grid joints, and NEWCOORD, where the object position coordinates are perturbed with a random offset. For RANDROB, all robots are placed in a connected manner, meaning that all objects can be reached by a robot.

BoxNet3D For BoxNet3D, we use MuJoCo to implement the feasibility checks such as robot arm collision and object collisions. The map size ranges from 2×2 to 4×4 , and the object number ranges from 1 to 4. For each map configuration, we randomly generate at most 100 different environments for the training data and at most 5 for the test data. Detailed dataset statistics are summarized in Table 6

B.3. A* Search Algorithm

We implement an A* search algorithm for solving the generated task. At each search step, the general workflow is: ① select the current best environment state, ② generate valid action for a single robot, ③ combine multiple valid actions and check whether they can run in parallel, and ④ update the environment with potential next step actions and put to candidate tools for next search step.

We list a Python reference code below:

Listing 1 | Reference A* search implementation

```
def astar_search(env: Any, max_iterations: int = 1000):
    open_set = []
    closed_set = set()

    g_scores: Dict[int, float] = {}
    came_from: Dict[int, Tuple[Optional[int], Optional[str]]] = {}
    states_cache: Dict[int, EnvStates] = {}

    try:
        initial_state_data = env.get_states()
        current_state = EnvStates(env, current_state_data=initial_state_data)
        initial_hash = current_state.hash()
    except Exception as e:
        raise e

    g_scores[initial_hash] = 0.0
    came_from[initial_hash] = (None, None)
    states_cache[initial_hash] = current_state

    heapq.heappush(open_set, (current_state.heuristic(), random.random(), initial_hash))
```

```

iterations = 0

while open_set and iterations < max_iterations:
    iterations += 1

    f_val, _, current_hash = heapq.heappop(open_set)

    if current_hash in closed_set:
        continue

    current_state = states_cache[current_hash]
    closed_set.add(current_hash)

    if current_state.is_goal():
        return reconstruct_path(came_from, current_hash)

    potential_next_moves = generate_potential_actions(current_state)

    for action_str, next_state_obj in potential_next_moves:
        if next_state_obj is None:
            continue

        next_hash = next_state_obj.hash()
        if next hash in closed_set:
            continue

        cost_of_this_action = 1.0
        tentative_g_score = g_scores[current_hash] + cost_of_this_action

        if tentative_g_score < g_scores.get(next_hash, float('inf')):
            came_from[next_hash] = (current_hash, action_str)
            g_scores[next_hash] = tentative_g_score
            states_cache[next_hash] = next_state_obj

        f_score_neighbor = tentative_g_score + next_state_obj.heuristic()
        heapq.heappush(open_set, (f_score_neighbor, random.random(), next_hash))

return None

class EnvStates:
    _hash_val: Optional[int] = None

    def __init__(self, env: Any, parent_state_data=None, current_state_data=None):
        self.env = env
        self.parent_state_data = parent_state_data
        self.cur_states = current_state_data
        # Assumes env has get_target_pos() and can define retraction height internally or via config
        self.target_positions = self.env.get_target_pos() # Target positions should include Z if
        ↪ relevant

    def hash(self) -> int:
        if self._hash_val is not None:
            return self._hash_val
        self._hash_val = xxhash.xxh64(self.cur_states.tobytes()).intdigest()
        return self._hash_val

    def _box_positions(self):
        self.env.reset(states=self.cur_states)
        return {
            boxname: self.env.get_box_pos(boxname)
            for boxname in self.env.object_names
        }

    def arm_positions(self):
        self.env.reset(states=self.cur_states)
        return {
            robot_name: self.env.get_arm_pos(robot_name)
            for robot_name in self.env.robot_names
        }

    def is_goal(self) -> bool:
        current_box_pos_map = self._box_positions()
        if len(current_box_pos_map) != len(self.target_positions):
            return False

        for box_name, target_val in self.target_positions.items():
            if box_name not in current_box_pos_map:
                return False
            # Use new env method for checking if object is at its target
            if not self.env.is_object_at_target(current_box_pos_map[box_name], target_val, box_name):

```

```

        return False
    return True

def heuristic(self) -> float:
    self.env.reset(states=self.cur_states)
    current_obj_positions_map = self._box_positions()

    h = 0.0
    num_matched = 0

    for name, target_pos_val in self.target_positions.items():
        if name in current_obj_positions_map:
            current_pos_val = current_obj_positions_map[name]
            if np.isnan(current_pos_val): # Check for NaN
                return float("inf")
            # Use env method for calculating distance/cost component for heuristic
            h += self.env.calculate_placement_quality(current_pos_val, target_pos_val, name)
            num_matched += 1
        else: # Object in target not found in current state
            return float("inf")

    if num_matched != len(self.target_positions): # Not all target objects were found or matched
        return float("inf")

    return math.sqrt(h) if h > 0 else 0.0

def apply_actions(self, action_strings: Union[List[str], str]) -> Optional["EnvStates"]:
    self.env.reset(states=self.cur_states)
    action_input = action_strings
    if isinstance(action_strings, list):
        action_input = "\n".join(action_strings)

    out = self.env.simulate_one_step(action_input)
    if out["success"]:
        return EnvStates(
            self.env,
            parent_state_data=self.cur_states,
            current_state_data=self.env.get_states(),
        )
    else:
        return None

# --- Utility Functions (Domain-specific helpers removed, env handles them) ---

def generate_single_robot_action(robot_id: str, state: EnvStates) -> List[str]:
    robot_actions = []
    env = state.env # Get the environment reference
    base_pos = env.get_base_pos(robot_id)
    arm_pos = state.arm_positions()[robot_id]

    for obj_id, obj_pos_val in state._box_positions().items():
        target_pos_val = state.target_positions[obj_id]

        if env.check_reach_range(robot_id, obj_pos_val): # Existing env call for reachability
            if env.is_object_at_target(obj_pos_val, target_pos_val, obj_id):
                continue

            # Get potential next positions for the object from the environment
            potential_next_obj_placements = env.get_valid_next_object_positions(
                obj_id, obj_pos_val, robot_id, base_pos
            )
            current_placement_quality = env.calculate_placement_quality(obj_pos_val, target_pos_val,
                obj_id
            )

            action_candidates_for_obj = []

            # Try to move the object to a better position
            for next_obj_p in potential_next_obj_placements:
                if env.calculate_placement_quality(next_obj_p, target_pos_val, obj_id) <
                    current_placement_quality:
                    # Check if arm is already at the object
                    if not env.is_arm_at_position(arm_pos, obj_pos_val[:2], robot_id):
                        # Action: Move arm to object
                        action_str = env.format_move_action_string(robot_id, obj_pos_val[:2], False)
                        action_candidates_for_obj.append(action_str)

            # Action: Move object (arm is now assumed to be at object or will be moved by
            # first action)
            action_str = env.format_move_action_string(robot_id, next_obj_p[:2], True)
            action_candidates_for_obj.append(action_str)

    return action_candidates_for_obj

```

```

# Try to move arm to an alternative/safe position if not productively moving an object
if not action_candidates_for_obj or all(
    # Check if any action involves carrying (True flag)
    # This logic might need refinement based on how format_move_action_string works
    # or if we have a better way to check if an action is "productive"
    "True" not in act for act in action_candidates_for_obj
):
    alternative_arm_destinations = env.get_alternative_arm_destinations(
        robot_id, base_pos, arm_pos
    )
    if alternative_arm_destinations:
        # Environment can decide which one to pick or return just one
        chosen_alt_dest = alternative_arm_destinations[0] # Take the first one
        if not env.is_arm_at_position(arm_pos, chosen_alt_dest[:2], robot_id):
            action_str = env.format_move_action_string(robot_id, chosen_alt_dest[:2],
                ↪ False)
            action_candidates_for_obj.append(action_str)

    for action_str_candidate in action_candidates_for_obj:
        if action_str_candidate not in robot_actions:
            robot_actions.append(action_str_candidate)

return robot_actions

def verify_parallel_actions(actionstr_input: Union[List[str], str], state: EnvStates) -> Tuple[bool,
→ Optional[EnvStates]]:
    current_env = state.env
    current_env.reset(states=state.cur_states)

    action_to_simulate = actionstr_input
    if isinstance(actionstr_input, list):
        action_to_simulate = "\n".join(actionstr_input)

    out = current_env.simulate_one_step(action_to_simulate)

    if out["success"]:
        new_state_data = current_env.get_states()
        new_search_state = EnvStates(
            current_env,
            parent_state_data=state.cur_states,
            current_state_data=new_state_data,
        )
        return True, new_search_state
    else:
        return False, None

def generate_potential_actions(state: EnvStates) -> List[Tuple[str, EnvStates]]:
    if not hasattr(state.env, 'robot_names'):
        return []

    robot_names = sorted(state.env.robot_names)

    single_robot_potential_actions: Dict[str, List[str]] = {
        r: generate_single_robot_action(r, state) for r in robot_names
    }

    valid_action_sets: List[Tuple[str, EnvStates]] = []
    action_verification_tasks = []

    for robot_id, actions in single_robot_potential_actions.items():
        for action in actions:
            action_verification_tasks.append(([action], state))

    active_robots = [r for r in robot_names if single_robot_potential_actions[r]]
    if len(active_robots) >= 2:
        max_concurrent_robots = 2
        if hasattr(state.env, 'object_names'): # Check if object_names exists before using its length
            max_concurrent_robots = min(4, len(state.env.object_names), len(active_robots))
        else: # Fallback if object_names is not available
            max_concurrent_robots = min(2, len(active_robots))

        for group_size in range(2, max_concurrent_robots + 1):
            if group_size > len(active_robots): continue
            for robot_group in itertools.combinations(active_robots, group_size):
                action_combos_for_group = list(
                    itertools.product(
                        *[single_robot_potential_actions[r] for r in robot_group]
                    )
    
```

```

        )
        cleaned_combos = [[a for a in combo] for combo in action_combos_for_group]
        for combo in cleaned_combos:
            action_verification_tasks.append((combo, state))

    verify_results = []
    for action_combo, original_state_for_verification in action_verification_tasks:
        success, new_state_obj = verify_parallel_actions(action_combo,
            ↪ original_state_for_verification)
        if success and new_state_obj is not None:
            verify_results.append(("\n".join(action_combo), new_state_obj))

    verify_results.sort(
        key=lambda x_tuple: (x_tuple[1].heuristic(), -x_tuple[0].count("True"),
            ↪ -len(x_tuple[0].split("\n"))))
    )

    return verify_results[:20]

def reconstruct_path(came_from: Dict[int, Tuple[Optional[int], str]], final_hash: int) -> List[str]:
    actions_sequence = []
    current_hash = final_hash
    while current_hash in came_from:
        parent_hash, action = came_from[current_hash]
        if action is not None:
            actions_sequence.append(action)
        if parent_hash is None:
            break
        current_hash = parent_hash
    actions_sequence.reverse()
    return actions_sequence

```

C. BoxNet2D and BoxNet3D Environment Examples

In this section, we provide more examples of the two environments we developed in this work. We note that all examples shown in this section are in 3×3 and 4×4 grid maps, but the dataset contains a wider range of map sizes. For more BoxNet3D video examples, please visit our project website at this anonymous link <https://github.com/UCSB-NLP-Chang/constraint-grounded-planner>.

C.1. BoxNet2D Examples

Original BoxNet2D examples Figure 6 shows two example BoxNet2D environments. Figure 7 shows four example collisions in BoxNet2D environments.

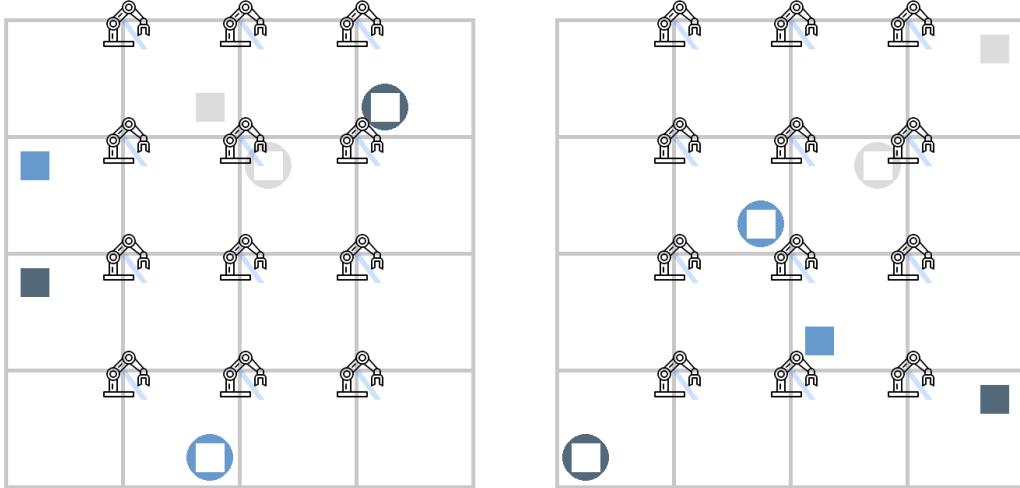


Figure 6 | Example BoxNet2D environment

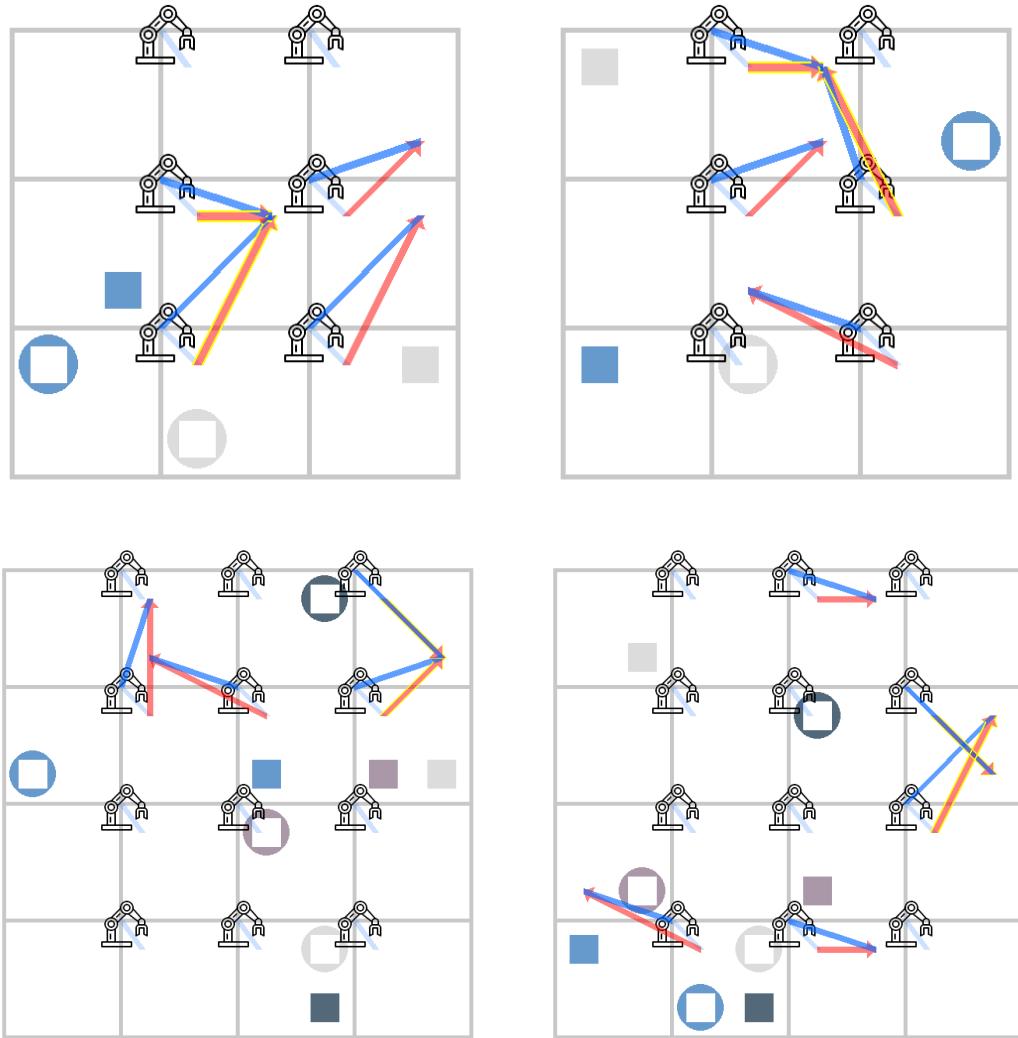


Figure 7 | Example collisions in BoxNet2D environment. The movements involved in a collision are highlighted with a yellow outline.

Unseen BoxNet2D examples for generalization experiment Figure 8 shows two example RANDROB environments. Figure 9 shows two example NEWCOORD environments.

C.2. BoxNet3D Examples

Figure 10 shows examples for BoxNet3D environments.

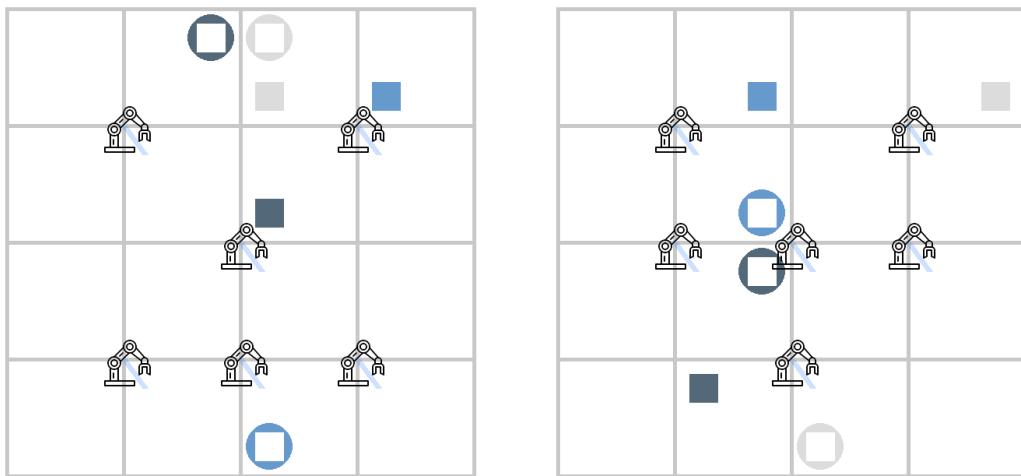


Figure 8 | Example RANDROB environment.

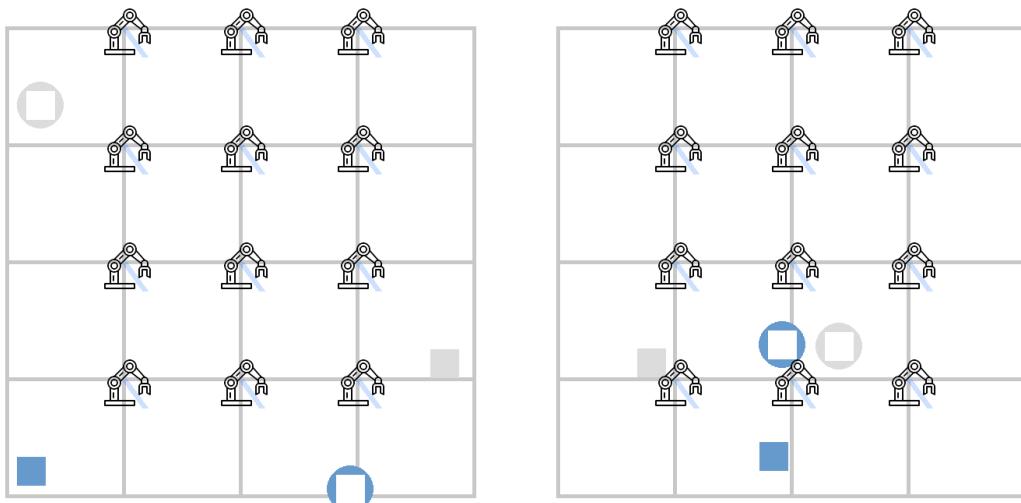


Figure 9 | Example NEWCOORD environment.

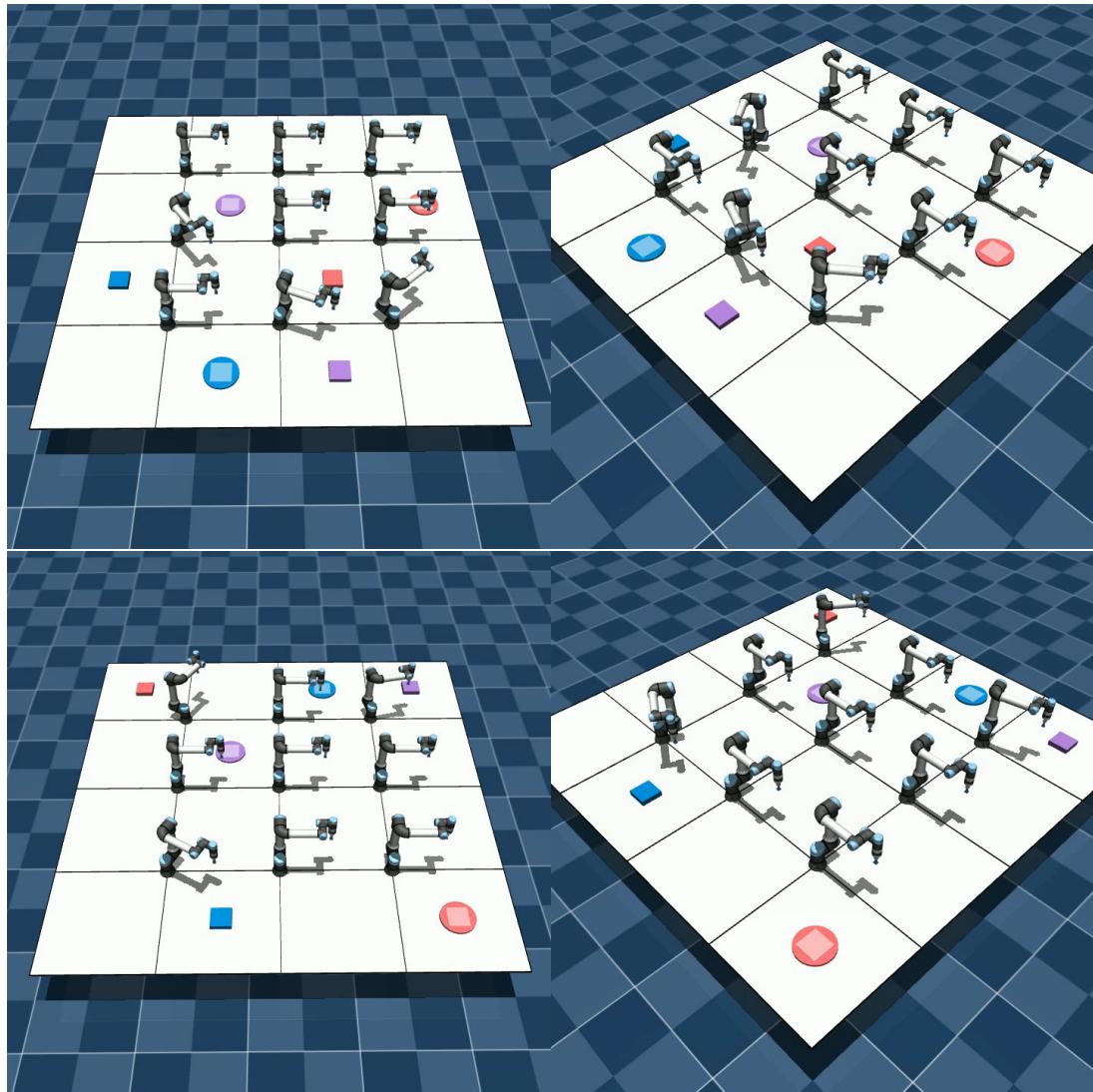


Figure 10 | Example BoxNet3D environment.

D. Detailed Prompts

In this section, we summarize the full prompts for SFT data synthesis in Section D.1, BoxNet2D and BoxNet3D environments with two planner modes in Section D.2, and reasoning behavior analysis in Section D.3.

D.1. Prompt for SFT Data Synthesis

We list the prompt for BoxNet2D thinking synthesis in Listing 2, and BoxNet3D thinking synthesis in Listing 3.

Listing 2 | Prompt for synthesizing reasoning trace for BoxNet2D

You are required to **assume the role of a central planner**. Your task is to simulate the
 ↳ step-by-step thinking process that logically leads you to the provided ground-truth plan.

Your thinking should be presented from a **first-person perspective**, clearly demonstrating your
 ↳ internal reasoning process of planning, validating and adjusting to avoid collision, and planning
 ↳ decisions.

Requirement for your generated first-person thinking:

1. **First-Person Perspective:** Write your internal thoughts as if you are personally making the
 ↳ decisions:
 - Use phrases like "Let me see...", "Wait, is that correct?", "I should check collisions
 ↳ first...", "Can I parallel two robot movements to make the plan more efficient?"
 - Demonstrate real-time analysis and potential hesitations or reconsiderations.
2. **Thinking Process with '<think>' Tags**:**
 - Enclose your entire reasoning sequence in `<think>` ... `</think>` tags.
 - Make sure you have explicit checks, e.g. collision checks, range feasibility, and confirmations
 ↳ of correctness. You can start the explicit checks with "Wait", "Hmm", "let me check", etc.
 - Make sure to pose questions to yourself, and then answer them. Show how you arrive at each
 ↳ movement decision.
 - You must include multiple explicit checks and self-questioning in your thinking process.

Below is the detailed task description. You can learn the rules for the task from these descriptions.

Task Description:

You are a central planner responsible for coordinating multiple robotic arms operating in a grid-like
 ↳ environment. Your goal is to plan and execute efficient, collision-free movements to transport
 ↳ objects to their designated target positions.

Task Representation:

- * Objective: Move all objects to their specified target locations safely and efficiently.
- * Input: A detailed map state containing positions of robots, objects, and target locations.
- * Output: A precise movement plan specifying each robot arm's actions for moving objects.

Position Representation:

- * All positions (robots, objects, targets) are given by their center coordinates, e.g., [0.25, 0.25],
 ↳ [0.75, 1.25].
- * Robots have a fixed base location and an extendable arm with a limited reach range.

Movement Rules:

- * Each robot arm can only move within a limited range relative to its fixed base position:
 - * X-axis: from (Base_X - 1.0) to (Base_X + 1.0) (exclusive).
 - * Y-axis: from (Base_Y - 1.0) to (Base_Y + 1.0) (exclusive).
- * For example:
 - * If a robot's base is [1.0, 1.0], its arm can reach [0.25, 0.75] or [1.25, 1.75], but not [0,
 ↳ 0.25] or [2.0, 0.75].
 - * Robots may move an object only if their arm aligns exactly with the object's current position,
 ↳ and if explicitly indicated in the action (move_object: True).

How to Generate Your Response:

Your response must **clearly indicate your thinking process** enclosed in `<think>` and `</think>` tags,
 ↳ followed by the generated step of the movement plan.

Thinking:

- * Clearly describe your analysis and decisions from a first-person perspective.
- * Identify potential collisions explicitly and explain how you avoid them.
- * Highlight your reasoning for movement choices, considering efficiency and collision avoidance.

Movement Plan (Output):

```

* Your generated step of the movement plan should be in markdown format and contain a JSON
  ↪ dictionary, with robot names as keys and their movement instructions as values, structured as
  ↪ follows:
```
json
[
{
 "robot_name": "start_position -> end_position, move_object"
 "robot_name": "start_position -> end_position, move_object",
},
{
 "robot_name": "start_position -> end_position, move_object",
}
]
```
* *start_position* and *end_position* represent the *[x, y]* coordinates of the robot arm's
  ↪ movement.
* *move_object* is a boolean indicating whether the robot moves an object (*True*) or simply moves
  ↪ its arm without carrying an object (*False*).
* Robots without actions in the current step should not be included.
Ensure your final step completes the objective of placing all objects at their target positions, and
  ↪ your plan forms a valid JSON array.

## Collision Avoidance Rules:
Your plan must strictly avoid collisions, as follows:
* Robot-Robot Collision:
  * Two robot arms cannot occupy the same position simultaneously.
  * Robot arms cannot intersect with each other or have intersecting movement trajectories
    ↪ during a step movement.
  * For example:
    * Collision occurs if Robot 1 moves [0.75, 0.75] -> 0.75, 1.25] and Robot 2 moves [2.25,
      ↪ 1.75] -> [0.75, 1.25] (same endpoint).
    * Collision occurs if Robot 1 moves [0.25, 0.25] -> [0.75, 0.25] and Robot 2 moves [1.25,
      ↪ 0.25] -> [0.25, 0.75] (intersecting arms as the end position Robot 1 is at the arm, as
      ↪ the end of Robot 2 arm position occupies [0.75, 0.25])
    * Collision occurs if Robot 1 moves [0.25, 0.25] -> [0.75, 0.75] and Robot 2 moves [0.25,
      ↪ 0.75] -> [0.75, 0.25] (intersecting movement as both arms moves across [0.5, 0.5]).
* Object-Object Collision:
  * Two objects cannot occupy the same position at any time.

## Example Environment and Ground-Truth Plan:
Below is an example scenario and its ground-truth solution:

```
text
{environment}
```

```

With the above information clearly provided, please start by explicitly presenting your first-person reasoning for the whole plan enclosed in <think>/</think> tags. Make sure you include explicit checks and self-questioning in your thinking process. Your reasoning should be clear and easy to follow, as if you are explaining it to someone else. Limit your thinking length within 2000 tokens.

Listing 3 | Prompt for synthesizing reasoning trace for BoxNet3D

You are required to **assume the role of a central planner**. Your task is to simulate the step-by-step thinking process that logically leads you to the provided ground-truth movement plan.

Your thinking should be presented from a **first-person perspective**, clearly demonstrating your internal reasoning process of planning, validating and adjusting to avoid collision, and planning decisions.

Requirement for your generated first-person thinking:

- **First-Person Perspective**:** Write your internal thoughts as if you are personally making the decisions:
 - Use phrases like "Let me see...", "Wait, is that correct?", "I should check collisions first...", "Can I parallel two robot movements to make the plan more efficient?"
 - Demonstrate real-time analysis and potential hesitations or reconsiderations.
- **Thinking Process with `<think>` Tags**:**
 - Enclose your entire reasoning sequence in `<think>` ... `</think>` tags.
 - Make sure you have explicit checks, e.g. collision checks, range feasibility, and confirmations of correctness. You can start the explicit checks with "Wait", "Hmm", "let me check", etc.
 - Make sure to pose questions to yourself, and then answer them. Show how you arrive at each movement decision.
 - You must include multiple explicit checks and self-questioning in your thinking process.

Below is the detailed task description. You can learn the rules for the task from these descriptions.

Task Description:

You are a central planner responsible for coordinating multiple robotic arms operating in a grid-like environment. Your goal is to plan and execute efficient, collision-free movements to transport objects to their designated target positions.

Task Representation:

- * Objective: Move all objects to their specified target locations safely and efficiently.
- * Input: A detailed map state containing positions of robots, objects, and target locations.
- * Output: A precise movement plan specifying each robot arm's actions for moving objects.

Position Representation:

- * All positions (robots, objects, targets) are given by their center coordinates, e.g., [0.55, 1.65], [2.75, 0.55].
- * Robots have a fixed base location and an extendable arm with a limited reach range.

Movement Rules:

- * Each robot arm can only move within a circular band around its fixed base position:
 - Let $d = \sqrt{(X - \text{Base}_X)^2 + (Y - \text{Base}_Y)^2}$.
 - The arm may reach (X, Y) only if $0.4 < d < 0.8$
- * For example:
 - If a robot's base is at [1.1, 1.1]:
 - It can reach [0.55, 0.55] since $\sqrt{(1.1 - 0.55)^2 + (1.1 - 0.55)^2} \approx 0.77 < 0.8$
 - It can reach [0.6, 1.1] since $\sqrt{(0.6 - 1.1)^2 + (1.1 - 1.1)^2} = 0.5 > 0.4$
 - It cannot reach [2.0, 1.1] because $\sqrt{(0.9^2 + 0^2)} = 0.9$, which exceeds 0.8
 - It cannot reach [2.25, 0.65] because $\sqrt{(1.15^2 + 0.45^2)} \approx 1.23$, which exceeds 0.8
 - If a robot needs to move an object within its range and the arm is not aligned with the object, the robot should first move its arm to the position of that object. By aligning, it measures the distance between object center and arm position is less than 0.1
 - When you plan a move, please follow following rules:
 - First check that the proposed target lies within the circular band $0.5 < d < 0.8$.
 - If it does not, adjust your plan or reject that movement.
 - If the arm is not yet aligned with an object it needs to move and that object lies within the band, plan a preliminary move to position the arm aligned with the object before picking it up.

How to Generate Your Response:

Your response must **clearly indicate your thinking process** enclosed in `<think>` and `</think>` tags, followed by the generated step of the movement plan.

Thinking:

- * Clearly describe your analysis and decisions from a first-person perspective.
- * Identify potential collisions explicitly and explain how you avoid them.
- * Highlight your reasoning for movement choices, considering efficiency and collision avoidance.

Movement Plan (Output):

- * Your generated step of the movement plan should be in markdown format and contain a JSON dictionary, with robot names as keys and their movement instructions as values, structured as follows:


```
```json
[
{
 "robot_name1": "Move end_position, move_object",
 "robot_name2": "Move end_position, move_object"
},
{
 "robot_name3": "Move end_position, move_object"
}
]
```
      * *end_position* represent the target [x, y] coordinates of the robot arm end point of the movement around circular path. Note that only the arm moves while its base remains fixed.
      * *move_object* is a boolean indicating whether the robot moves an object (*True*) or simply moves its arm without carrying an object (*False*).
      * One robot can only be moved once in each step, which means that no repeated keys are allowed in the same step.
      * Robots without actions in the current step should not be included.
      Ensure your final step completes the objective of placing all objects at their target positions, and your plan forms a valid JSON array.
    
```

Collision Avoidance Rules:

Your plan must strictly avoid collisions, as follows:

- * Robot-Robot Collision

```

* Each robot arm always swings along a smooth **circular** path around its base.
* Two robot arms cannot occupy the same position at the end of a move.
* Their curved paths must not cross or share any point during the move.
* Sometimes a robot needs to move its arm to a safe position to avoid collision between another
  ↳ robot that move its arm to reach an object.
* Example:
  * robot_0 swings from [0.25, 0.25] to [0.75, 0.75] and robot_1 swings from [0.25, 0.75] to [0.75,
    ↳ 0.25] at the same time. Both arcs pass through [0.5, 0.5], causing a collision.
* Object-Object Collision
  * Two objects cannot occupy the same (x, y) at any time.
  * If you move more than one object at once, they must have different drop-off points and
    ↳ non-crossing straight-line paths.
* Robot-Object Collision
  * An arm's circular path must not sweep through any object it isn't carrying.
  * Before moving, confirm the curved trajectory does not pass over another object's position.

## Plan Efficiency Considerations:
* Each step of your plan involves simultaneous robot arm movements from their current positions to
  ↳ specific target positions.
* Each robot arm moves at a constant speed of 0.5 units/time.
* The duration of each step is determined by the longest single-arm movement within that step.
* The total execution time is the sum of all individual step durations.
* You should aim to minimize total execution time while ensuring collision-free movements and
  ↳ successful object placements.

## Example Environment and Ground-Truth Plan:
Below is an example scenario and its ground-truth solution:

```text
{environment}
```

```

With the above information clearly provided, please start by explicitly presenting your first-person
 ↳ reasoning for the whole plan enclosed in <think></think> tags. Make sure you include explicit
 ↳ checks and self-questioning in your thinking process. Your reasoning should be clear and easy to
 ↳ follow, as if you are explaining it to someone else. Limit your thinking length within 2000
 ↳ tokens.

D.2. Prompt for BoxNet2D and BoxNet3D environment

BoxNet2D We list the prompt for FULLPLAN planner in BoxNet2D in List 4, and REPLAN planner in List 5.

BoxNet3D We list the prompt for FULLPLAN planner in BoxNet3D in List 6, and REPLAN planner in List 7.

Listing 4 | Prompt for FULLPLAN planner in BoxNet2D environment

You are a central planner responsible for coordinating robotic arms in a grid-like environment to
 ↳ transport objects to their designated targets. Each robot is stationed at the corner of a 1x1
 ↳ square and uses its arm to move objects. Your task is to generate an efficient and collision-free
 ↳ plan for multiple robots, ensuring all objects reach their target positions after the whole plan
 ↳ is executed.

```

## Task Description:
*Task Representation:
* Objective: Move all objects to their specified target locations safely and efficiently.
* Input: A detailed map state containing positions of robots, objects, and target locations.
* Output: A precise movement plan specifying each robot arm's actions for moving objects.

*Position Representation:
* All positions (robots, objects, targets) are given by their center coordinates, e.g., [0.25, 0.25],
  ↳ [0.75, 1.25].
* Robots have a fixed base location and an extendable arm with a limited reach range.

*Movement Rules:
Your generated movement must strictly consider the reachability of each robot arm, detaild rule in
  ↳ following:
* Each robot arm can only move within a limited range relative to its fixed base position:
  * X-axis: from (Base_X - 1.0) to (Base_X + 1.0) (exclusive).

```

```

    * Y-axis: from (Base_Y - 1.0) to (Base_Y + 1.0) (exclusive).
* For example:
  * If a robot's base is [1.0, 1.0], its arm can reach [0.25, 0.75] or [1.25, 1.75], but not [0, → 0.25] or [2.25, 1.75] because  $2.25 - 1.0 = 1.25 > 1.0$  and  $0 - 1.0 = -1.0 \leq -1$ .
  * Robots may move an object only if their arm position aligns exactly with the object's current → position, and if explicitly indicated in the action (move_object: True).
  * Make sure you explicitly think about whether your proposed movement for one arm is valid, and → correct it if it is not.
  * If a robot needs to move an object within its range and the arm is not aligned with the object, → the robot should first move its arm to the position of that object.

## How to Generate Your Response:
Your response must **clearly indicate your thinking process** enclosed in <think> and </think> tags, → followed by the generated step of the movement plan.

*Thinking:
* Clearly describe your analysis and decisions from a first-person perspective.
* You should think carefully whether your plan has collision by explicitly generating your thoughts, → and avoid them in your final output if there is any.
* Highlight your reasoning for movement choices, considering efficiency and collision avoidance.

*Movement Plan (Output):
* Your generated step of the movement plan should be in markdown format and contain a JSON list, with → each entry as a dictionary indicating one step, and the robot names are keys and their movement → instructions as values for each step, structured as follows:
```json
[
{
 "robot_name1": "start_position -> end_position, move_object",
 "robot_name2": "start_position -> end_position, move_object"
},
{
 "robot_name3": "start_position -> end_position, move_object"
}
]
```
* *start_position* and *end_position* represent the *[x, y]* coordinates of the robot arm before and → after the movement. Note that only the arm moves while base remains fixed.
* *move_object* is a boolean indicating whether the robot moves an object (*True*) or simply moves → its arm without carrying an object (*False*).
* Robots without actions in a certain step should not be included.
* One robot can only be moved once in each step, which means that no repeated keys are allowed in the → same step.
Ensure your final step completes the objective of placing all objects at their target positions, and → your plan forms a valid JSON list.

## Collision Avoidance Rules:
Your plan must strictly avoid collisions, as follows:
* Robot-Robot Collision:
  * Two robot arms cannot occupy the same position simultaneously.
  * Robot arms cannot intersect with each other or have intersecting movement trajectories → during a step movement.
  * For example:
    * Collision occurs if Robot 1 moves [0.75, 0.75] -> 0.75, 1.25] and Robot 2 moves [2.25, → 1.75] -> [0.75, 1.25] (same endpoint).
    * Collision occurs if Robot 1 moves [0.25, 0.25] -> [0.75, 0.25] and Robot 2 moves [1.25, → 0.25] -> [0.25, 0.75] (intersecting arms as the end position Robot 1 is at the arm, as → the end of Robot 2 arm position occupies [0.75, 0.25])
    * Collision occurs if Robot 1 moves [0.25, 0.25] -> [0.75, 0.75] and Robot 2 moves [0.25, → 0.75] -> [0.75, 0.25] (intersecting movement as both arms moves across [0.5, 0.5]).
  * Object-Object Collision:
    * Two objects cannot occupy the same position at any time.

## Plan Efficiency Considerations:
* Each step of your plan involves simultaneous robot arm movements from their current positions to → specified target positions.
* Each robot arm moves at a constant speed of 0.5 units/time.
* The duration of each step is determined by the longest single-arm movement within that step.
* The total execution time is the sum of all individual step durations.
* You should aim to minimize total execution time while ensuring collision-free movements and → successful object placements.

## Example Input & Output:
* Input:
Object positions:
  Object 1: [0.75, 0.75]
```

```

Object 2: [1.75, 0.25]
Target positions:
  Object 1 target: [2.25, 0.75]
  Object 2 target: [0.25, 1.25]
Robot positions:
  Robot 1: base [1.0, 1.0], arm [0.75, 0.75]
  Robot 2: base [2.0, 0.0], arm [1.75, 0.75]

* Output:
<think> Let's understand the scenerio ... </think>
```
[{"Robot 1": "[0.75, 0.75] -> [1.25, 0.25], True", "Robot 2": "[1.75, 0.75] -> [1.75, 0.25], False"}, {"Robot 2": "[1.75, 0.25] -> [1.75, 0.75], True"}, {"Robot 1": "[1.25, 0.25] -> [1.75, 0.75], False", "Robot 2": "[1.75, 0.75] -> [1.25, 0.25], False"}, {"Robot 1": "[1.75, 0.75] -> [0.25, 1.25], True", "Robot 2": "[1.25, 0.25] -> [2.25, 0.75], True"}]
```

Given the information above, now consider the following environment:
Input:

{mapstate}

Generate the full plan for moving these robots.

```

Listing 5 | Prompt for REPLAN planner in BoxNet2D environment

You are a central planner responsible for coordinating robotic arms in a grid-like environment to transport objects to their designated targets. Each robot is stationed at the corner of a 1x1 square and uses its arm to move objects. Your task is to interactively generate an efficient and collision-free movement plan for controlling these robots, targeting at moving all objects to their target positions.

At each step, you will receive the current state of the environment wrapped by `<observation>` and `</observation>` tags. You need to generate the next-step plan for moving the robots, ensuring that your output contains your thinking process and the markdown json dict.

```

## Task Description:
*Task Representation:*
* Objective: Move all objects to their specified target locations safely and efficiently.
* Input: A detailed map state containing positions of robots, objects, and target locations.
* Output: A precise movement plan specifying each robot arm's actions for moving objects.

*Position Representation:*
* All positions (robots, objects, targets) are given by their center coordinates, e.g., [0.25, 0.25],
→ [0.75, 1.25].
* Robots have a fixed base location and an extendable arm with a limited reach range.

*Movement Rules:*
* Each robot arm can only move within a limited range relative to its fixed base position:
  * X-axis: from (Base_X - 1.0) to (Base_X + 1.0) (exclusive).
  * Y-axis: from (Base_Y - 1.0) to (Base_Y + 1.0) (exclusive).
* For example:
  * If a robot's base is [1.0, 1.0], its arm can reach [0.25, 0.75] or [1.25, 1.75], but not [0,
→ 0.25] or [2.25, 1.75] because 2.25 - 1.0 = 1.25 > 1.0 and 0 - 1.0 = -1.0 <= -1.
  * Robots may move an object only if their arm aligns exactly with the object's current position,
→ and if explicitly indicated in the action (move_object: True).
  * If a robot needs to move an object within its range and the arm is not aligned with the object,
→ the robot should first move its arm to the position of that object.

## How to Generate Your Response:
Your response must **clearly indicate your thinking process** enclosed in <think> and </think> tags,
→ followed by the generated step of the movement plan.

*Thinking:*
* Clearly describe your analysis and decisions from a first-person perspective.
* Think carefully and try to identify potential collisions explicitly in the analysis and explain how
→ you avoid them.
* Highlight your reasoning for movement choices, considering efficiency and collision avoidance.

```

```

*Movement Plan (Output):*
* Your generated step of the movement plan should be in markdown format and contain a JSON
→ dictionary, with robot names as keys and their movement instructions as values, structured as
→ follows:
```json
{
 "robot_name1": "start_position -> end_position, move_object"
 "robot_name2": "start_position -> end_position, move_object",
},
```
* *start_position* and *end_position* represent the *[x, y]* coordinates of the robot arm before and
→ after the movement. Note that only the arm moves while base remains fixed.
* *move_object* is a boolean indicating whether the robot moves an object (*True*) or simply moves
→ its arm without carrying an object (*False*).
* Robots without actions in a certain step should not be included.
* One robot can only be moved once in each step, which means that no repeated keys are allowed in the
→ same step.
Ensure your output forms a valid JSON dictionary of next-step plan.

## Collision Avoidance Rules:
Your plan must strictly avoid collisions, as follows:
* Robot-Robot Collision:
    * Two robot arms cannot occupy the same position simultaneously.
    * Robot arms cannot intersect with each other or have intersecting movement trajectories
→ during a step movement.
    * For example:
        * Collision occurs if Robot 1 moves [0.75, 0.75] -> 0.75, 1.25] and Robot 2 moves [2.25,
→ 1.75] -> [0.75, 1.25] (same endpoint).
        * Collision occurs if Robot 1 moves [0.25, 0.25] -> [0.75, 0.25] and Robot 2 moves [1.25,
→ 0.25] -> [0.25, 0.75] (intersecting arms as the end position Robot 1 is at the arm, as
→ the end of Robot 2 arm position occupies [0.75, 0.25])
        * Collision occurs if Robot 1 moves [0.25, 0.25] -> [0.75, 0.75] and Robot 2 moves [0.25,
→ 0.75] -> [0.75, 0.25] (intersecting movement as both arms moves across [0.5, 0.5]).
* Object-Object Collision:
    * Two objects cannot occupy the same position at any time.

## Plan Efficiency Considerations:
* The execution time of your plan involves simultaneous robot arm movements from their current
→ positions to specified target positions.
* Each robot arm moves at a constant speed of 0.5 units/time.
* The duration of the plan is determined by the longest single-arm movement within it.
* You should aim to minimize the execution time while ensuring collision-free movements and
→ successful object placements.

## Example Input & Output:
Input:
<observation>
Object positions:
    Object 1: [0.75, 0.75]
    Object 2: [1.75, 0.25]
Target positions:
    Object 1 target: [2.25, 0.75]
    Object 2 target: [0.25, 1.25]
Robot positions:
    Robot 1: base [1.0, 1.0], arm [0.75, 0.75]
    Robot 2: base [2.0, 0.0], arm [1.75, 0.75]
</observation>

* Output:
<think> Let's understand the scenario ... </think>
```json
{"Robot 1": "[0.75, 0.75] -> [1.25, 0.25], True", "Robot 2": "[1.75, 0.75] -> [1.75, 0.25], False"}
```

Now work on the following problem given by user.

<observation>
{mapstate}
</observation>

```

Listing 6 | Prompt for FULLPLAN planner in BoxNet3D environment

You are a central planner responsible for coordinating robotic arms in a grid-like environment to
 → transport objects to their designated targets. Each robot is stationed at the corner of a 1.1x1.1
 → square and uses its arm to move objects. Your task is to generate an efficient and collision-free
 → plan for multiple robots, ensuring all objects reach their target positions after the whole plan
 → is executed.

```
## Task Description:  

*Task Representation:  

* Objective: Move all objects to their specified target locations safely and efficiently.  

* Input: A detailed map state containing positions of robots, objects, and target locations.  

* Output: A precise movement plan specifying each robot arm's actions for moving objects.

*Position Representation:  

* All positions (robots, objects, targets) are given by their center coordinates, e.g., [0.55, 1.65],  

  → [2.75, 0.55].  

* Robots have a fixed base location and an extendable arm with a limited reach range.

*Movement Rules:  

Your generated movement must strictly consider the reachability of each robot arm, detaild rule in  

→ following:  

* Each robot arm can only move within a circular band around its fixed base position:  

  - Let  $d = \sqrt{((X - \text{Base}_X)^2 + (Y - \text{Base}_Y)^2)}$ .  

  - The arm may reach  $(X, Y)$  only if  $0.4 < d < 0.8$ 

* For example:  

  - If a robot's base is at [1.1, 1.1]:  

    - It can reach [0.55, 0.55] since  $\sqrt{(1.1 - 0.55)^2 + (1.1 - 0.55)^2} \approx 0.77 < 0.8$   

    - It can reach [0.6, 1.1] since  $\sqrt{(0.6 - 1.1)^2 + (1.1 - 1.1)^2} = 0.5 > 0.4$   

    - It cannot reach [2.0, 1.1] because  $\sqrt{(2.0 - 1.1)^2 + 0^2} = 0.9$ , which exceeds 0.8  

    - It cannot reach [2.25, 0.65] because  $\sqrt{(2.25 - 1.1)^2 + (0.65 - 1.1)^2} \approx 1.23$ , which exceeds 0.8  

  - If a robot needs to move an object within its range and the arm is not aligned with the object,  

  → the robot should first move its arm to the position of that object. By aligning, it meas the  

  → distance between object center and arm position is less than 0.1  

  - When you plan a move, please follow following rules:  

    - First check that the proposed target lies within the circular band  $0.5 < d < 0.8$ .  

    - If it does not, adjust your plan or reject that movement.  

    - If the arm is not yet aligned with an object it needs to move and that object lies within the  

  → band, plan a preliminary move to position the arm aligned with the object before picking it  

  → up.

## How to Generate Your Response:  

Your response must **clearly indicate your thinking process** enclosed in <think> and </think> tags,  

→ followed by the generated step of the movement plan.

*Thinking:  

* Clearly describe your analysis and decisions from a first-person perspective.  

* You should think carefully whether your plan has collision by explicitly generating your thoughts,  

  → and avoid them in your final output if there is any.  

* Highlight your reasoning for movement choices, considering efficiency and collision avoidance.

*Movement Plan (Output):  

* Your generated step of the movement plan should be in markdown format and contain a JSON list, with  

  → each entry as a dictionary indicating one step, and the robot names are keys and their movement  

  → instructions as values for each step, structured as follows:  

```json
[

{
 "robot_name1": "Move end_position, move_object",

 "robot_name2": "Move end_position, move_object"
},

{
 "robot_name3": "Move end_position, move_object"
}

]
```
* *end_position* represent the target *[x, y]* coordinates of the robot arm after the movement. Note  

  → that only the arm moves while its base remains fixed.  

* *move_object* is a boolean indicating whether the robot moves an object (*True*) or simply moves  

  → its arm without carrying an object (*False*).  

* Robots without actions in a certain step should not be included.  

* One robot can only be moved once in each step, which means that no repeated keys are allowed in the  

  → same step.  

Ensure your final step completes the objective of placing all objects at their target positions, and  

→ your plan forms a valid JSON list.
```

```

## Collision Avoidance Rules:
Your plan must strictly avoid collisions, as follows:
* Robot-Robot Collision
  * Each robot arm always swings along a smooth **circular** path around its base.
  * Two robot arms cannot occupy the same position at the end of a move.
  * Their curved paths must not cross or share any point during the move.
  * Sometimes a robot needs to move its arm to a safe position to avoid collision between another
    → robot that moves its arm to reach an object.
* **Example:** 
  * robot_0 swings from [0.25, 0.25] to [0.75, 0.75] and robot_1 swings from [0.25, 0.75] to [0.75,
    → 0.25] at the same time. Both arcs pass through [0.5, 0.5], causing a collision.
* Object-Object Collision
  * Two objects cannot occupy the same (x, y) at any time.
  * If you move more than one object at once, they must have different drop-off points and
    → non-crossing straight-line paths.
* Robot-Object Collision
  * An arm's circular path must not sweep through any object it isn't carrying.
  * Before moving, confirm the curved trajectory does not pass over another object's position.

## Plan Efficiency Considerations:
* Each step of your plan involves simultaneous robot arm movements from their current positions to
  → specified target positions.
* Each robot arm moves at a constant speed of 0.5 units/time.
* The duration of each step is determined by the longest single-arm movement within that step.
* The total execution time is the sum of all individual step durations.
* You should aim to minimize total execution time while ensuring collision-free movements and
  → successful object placements.

## Example Input & Output:
* Input:
Object positions:
  Object 1: [0.55, 1.65]
Target positions:
  Object 1 target: [1.65, 0.55]
Robot positions:
  Robot 1: base [1.1, 1.1], arm [1.24, 0.61]
  Robot 2: base [1.1, 2.2], arm [1.24, 1.71]

* Output:
<think> Let's understand the scenario ... </think>
```
json
[
 {
 "Robot 1": "Move [0.55, 1.65] False"
 },
 {
 "Robot 1": "Move [1.65, 1.65] True"
 },
 {
 "Robot 1": "Move [1.10, 1.70] False", "Robot 0": "Move [1.65, 1.66] False"
 },
 {
 "Robot 0": "Move [1.65, 0.55] True"
 }
]
```

Given the information above, now consider the following environment:
Input:
{mapstate}
Generate the full plan for moving these robots.

```

Listing 7 | Prompt for REPLAN planner in BoxNet3D environment

You are a central planner responsible for coordinating robotic arms in a grid-like environment to
 → transport objects to their designated targets. Each robot is stationed at the corner of a 1.1x1.1
 → square and uses its arm to move objects. Your task is to interactively generate an efficient and
 → collision-free movement plan for controlling these robots, targeting at moving all objects to
 → their target positions.

At each step, you will receive the current state of the environment wrapped by <observation> and
 → </observation> tags. You need to generate the next-step plan for moving the robots, ensuring that
 → your output contains your thinking process and the markdown json dict.

```

## Task Description:
*Task Representation:*
* Objective: Move all objects to their specified target locations safely and efficiently.
* Input: A detailed map state containing positions of robots, objects, and target locations.
* Output: A precise movement plan specifying each robot arm's actions for moving objects.

*Position Representation:*
* All positions (robots, objects, targets) are given by their center coordinates, e.g., [0.55, 1.65],
→ [2.75, 0.55].
* Robots have a fixed base location and an extendable arm with a limited reach range.

*Movement Rules:*
Your generated movement must strictly consider the reachability of each robot arm, detailed rule in
→ following:
* Each robot arm can only move within a circular band around its fixed base position:
- Let  $d = \sqrt{(X - \text{Base}_X)^2 + (Y - \text{Base}_Y)^2}$ .
- The arm may reach  $(X, Y)$  only if  $0.4 < d < 0.8$ 

* For example:
- If a robot's base is at [1.1, 1.1]:
  - It can reach [0.55, 0.55] since  $\sqrt{(1.1 - 0.55)^2 + (1.1 - 0.55)^2} \approx 0.77 < 0.8$ 
  - It can reach [0.6, 1.1] since  $\sqrt{(0.6 - 1.1)^2 + (1.1 - 1.1)^2} = 0.5 > 0.4$ 
  - It cannot reach [2.0, 1.1] because  $\sqrt{(2.0 - 1.1)^2 + (1.1 - 1.1)^2} = 0.9$ , which exceeds 0.8
  - It cannot reach [2.25, 0.65] because  $\sqrt{(2.25 - 1.1)^2 + (0.65 - 1.1)^2} \approx 1.23$ , which exceeds 0.8
- If a robot needs to move an object within its range and the arm is not aligned with the object,
→ the robot should first move its arm to the position of that object. By aligning, it means the
→ distance between object center and arm position is less than 0.1
- When you plan a move, please follow following rules:
  - First check that the proposed target lies within the circular band  $0.4 < d < 0.8$ .
  - If it does not, adjust your plan or reject that movement.
  - If the arm is not yet aligned with an object it needs to move and that object lies within the
→ band, plan a preliminary move to position the arm aligned with the object before picking it
→ up.

## How to Generate Your Response:
Your response must **clearly indicate your thinking process** enclosed in <think> and </think> tags,
→ followed by the generated step of the movement plan.

*Thinking:*
* Clearly describe your analysis and decisions from a first-person perspective.
* Think carefully and try to identify potential collisions explicitly in the analysis and explain how
→ you avoid them.
* Highlight your reasoning for movement choices, considering efficiency and collision avoidance.

*Movement Plan (Output):*
* Your generated step of the movement plan should be in markdown format and contain a JSON
→ dictionary, with robot names as keys and their movement instructions as values, structured as
→ follows:
```json
{
 "robot_name1": "Move end_position, move_object",
 "robot_name2": "Move end_position, move_object"
},
```
* *end_position* represent the target *[x, y]* coordinates of the robot arm end point of the
→ movement around circular path. Note that only the arm moves while its base remains fixed.
* *move_object* is a boolean indicating whether the robot moves an object (*True*) or simply moves
→ its arm without carrying an object (*False*).
* Robots without actions in a certain step should not be included.
* One robot can only be moved once in each step, which means that no repeated keys are allowed in the
→ same step.
Ensure your output forms a valid JSON dictionary of next-step plan.

## Collision Avoidance Rules:
Your plan must strictly avoid collisions, as follows:
* Robot-Robot Collision
  * Each robot arm always swings along a smooth **circular** path around its base.
  * Two robot arms cannot occupy the same position at the end of a move.
  * Their curved paths must not cross or share any point during the move.
  * Sometimes a robot needs to move its arm to a safe position to avoid collision between another
→ robot that moves its arm to reach an object.
* ***Example:**
```

```

* robot_0 swings from [0.25, 0.25] to [0.75, 0.75] and robot_1 swings from [0.25, 0.75] to [0.75,
→ 0.25] at the same time. Both arcs pass through [0.5, 0.5], causing a collision.
* Object-Object Collision
  * Two objects cannot occupy the same (x, y) at any time.
  * If you move more than one object at once, they must have different drop-off points and
→ non-crossing straight-line paths.
* Robot-Object Collision
  * An arm's circular path must not sweep through any object it isn't carrying.
  * Before moving, confirm the curved trajectory does not pass over another object's position.

## Plan Efficiency Considerations:
* The execution time of your plan involves simultaneous robot arm movements from their current
→ positions to specified target positions.
* Each robot arm moves at a constant speed of 0.5 units/time.
* The duration of the plan is determined by the longest single-arm movement within it.
* You should aim to minimize the execution time while ensuring collision-free movements and
→ successful object placements.

## Example Input & Output:
Input:
<observation>
Object positions:
  Object 1: [0.55, 1.65]
Target positions:
  Object 1 target: [1.65, 0.55]
Robot positions:
  Robot 1: base [1.1, 1.1], arm [1.24, 0.61]
  Robot 2: base [1.1, 2.2], arm [1.24, 1.71]
</observation>

* Output:
<think> Let's understand the scenario ... </think>
```json
{
 "Robot 1": "Move [0.55, 1.65] False"
}
```

Now work on the following problem given by user:

<observation>
{mapstate}
</observation>

```

D.3. Prompt for Reasoning Behavior Probing

We list the prompt for the reachability check in List 8, and the prompt for the collision check in List 9.

Listing 8 | Prompt for GPT-4o to count reachability check

How many reachability checks about the robot's movement are presented in the following reasoning trace? For example, a sentence like 'Robot 0 (base [1.0, 1.0]) can reach [0.25, 0.25]' counts as one verification. Give me an integer number without saying anything else.
The reasoning trace is:
{trace}

Listing 9 | Prompt for GPT-4o to count collision check

How many collision checks about the robot's movement are presented in the following reasoning trace?
For example, a sentence like 'Robot 0 moves to [0.25, 0.25], and Robot 1 moves to [0.25, 0.25]. They may collide with each other.' counts as one verification. Give me an integer number without saying anything else.
The reasoning trace is:
{trace}