BAR: A Backward Reasoning based Agent for Complex Minecraft Tasks

Weihong Du¹² Wenrui Liao¹² Binyu Yan^{1*} Hongru Liang¹² Anthony G. Cohn³⁴ Wenqiang Lei¹² ¹College of Computer Science, Sichuan University, China ²Engineering Research Center of Machine Learning and Industry Intelligence, Ministry of Education, China ³School of Computer Science, University of Leeds, UK ⁴The Alan Turing Institute, UK {duweihong, liaowenrui}@stu.scu.edu.cn {yby, lianghongru, wengianglei}@scu.edu.cn

riangnongiu, wenqiangiei;escu.eo

a.g.cohn@leeds.ac.uk

Abstract

Large language model (LLM) based agents have shown great potential in following human instructions and automatically completing various tasks; to do this, the agent needs to decompose it into easily executed steps by planning. Existing LLM-based approaches to planning mostly proceed by inferring what steps should be inserted into the plan next by starting from the agent's initial state. However, this forward reasoning paradigm does not work well for complex tasks. We study this issue in Minecraft, a virtual environment that simulates complex tasks based on real-world scenarios. The failure of forward reasoning is often caused by the large perception gap between the agent's initial state and task goal. To alleviate this, we leverage backward reasoning and make the planning start from the terminal (or goal) state, by first considering which actions could directly achieve the task goal in one step, before proceeding to consider how the preconditions of those actions can in turn be achieved. Our BAckward Reasoning based agent (BAR) is equipped with a recursive goal decomposition module, a state consistency maintaining module and a stage memory module. Experimental results demonstrate the superiority of **BAR** over existing methods and the effectiveness of proposed modules. The code and dataset are available in https://github.com/SCUNLP/BAR.

1 Introduction

Large language model (LLM) based agents are a powerful tool to support human-machine intelligence because of their potential to follow human instructions. However, it is not possible to complete one complex task in a single step. The key lies in how to decompose the task into easily executed steps. Existing LLM-based approaches mostly focus on forward reasoning based planning (Huang et al., 2024; Koh et al., 2024; Yu et al., 2024a; Chen et al., 2024). They infer what steps should be executed next starting from the agent's initial state, where no steps have yet been executed. However, we observe this reasoning paradigm doesn't work well for complex tasks. Let us take a close look at a complex task (obtain a diamond pickaxe) in an open-world environment, Minecraft, which simulates numerous complex tasks based on real-world scenarios. As shown in the upper part of Figure 1, the forward reasoning agent fails to infer the correct next step after the step "Craft 8 / (stick)". This is because there is a huge perception gap between the agent's initial state and the task goal.

To tackle this issue, we propose to use backward reasoning (Yu et al., 2024b), another powerful reasoning paradigm¹. As shown in the lower part of Figure 1, by backward reasoning, the agent infers the steps to be executed starting from the terminal state. In the terminal state, the task goal has been achieved and the agent has already obtained 1 \nearrow (diamond pickaxe). To achieve this goal, it is easy to infer that the last step should be "Craft 1 $\cancel{3}$ " because after executing this step the agent can obtain 1 diamond pickaxe. Furthermore, to successfully execute the step "Craft 1 $\cancel{3}$ ", the agent needs to prepare the materials for it. So the second to last step should be "Mine 3 (diamond)" as (a) is the

¹In fact, planning by reasoning backwards has a long history outside of the LLM community, dating back to Newell et al. (1959), who articulated their "principle of subgoal reduction" in their General Problem Solver (GPS); subsequently "regression planning" became a major focus in the symbolic planning community (e.g. (McDermott, 1991)).

^{*} Corresponding author



Figure 1: A toy example showing the essential advantage of backward reasoning over forward reasoning for task planning. In Minecraft, most tasks are started by "Mine 🗣 (log)". Best viewed in color.



Figure 2: The conflict of the agent's physical states during planning. In Minecraft, (stone) can only be obtained below the ground, so it is necessary to reach below the ground first and then mine (stone). Best viewed on screen.

key material for crafting a diamond pickaxe. Subsequently, the agent can continue to infer the third to last step to support the execution of step "Mine 3 ④", repeating until the entire plan is generated.

Despite this advantage, backward reasoning is not a perfect solution — the agent only considers the current goal to be achieved, ignoring the current world state. As shown in Figure 2, the agent infers the step "Mine (stone)" immediately after inserting the step "Craft 1 \nearrow (wooden pickaxe)". However, this step cannot be executed successfully because (stone) can only be obtained below the ground but the agent was above the ground before executing this step. At last, this conflict results in a failing planning for the task.

To this end, we propose **BAR**, a <u>BA</u>ckward <u>R</u>easoning based agent for complex Minecraft tasks, which is equipped with three carefullydesigned modules. (1) Inspired by the dynamic programming algorithm (Bellman, 1954), we design the recursive goal decomposition module to support robust backward planning starting from the terminal state. (2) We develop the state consistency maintaining module to monitor the agent's states and handle any state conflicts during planning. (3) We introduce the stage memory module to make the best of existing interactions and promote more efficient planning. Experimental results show that the **BAR** agent outperforms other SOTA methods by a large margin. We hope our exploration and discoveries can promote the development of more efficient LLM-based planning agents. In summary, our contributions are three-fold:

- We emphasize backward reasoning over forward reasoning for complex tasks, as it can avoid the issue caused by the perception gap between the agent's initial state and the task goal.
- We propose **BAR** and make robust, consistent, and efficient planning from the recursive goal decomposition, state consistency maintaining, and stage memory modules.
- Experimental results in Minecraft demonstrate the superiority of **BAR** and the effectiveness of the proposed modules.

2 Related Work

2.1 Task Planning by Agent

Task planning aims at decomposing tasks into steps which can be directly executed. Current LLMbased approaches mostly focus on forward reasoning based planning. Some studies conduct iterative forward reasoning to generate the steps to be executed starting from the agent's initial state (Sanyal et al., 2022; Huang et al., 2022; Xue et al., 2023; Jiang et al., 2024; Koh et al., 2024; Yu et al., 2024a). However, these studies fail to generate correct plans due to the big perception gap between the agent's initial state and task goal. To reduce the difficulty of planning, some studies choose the best steps from a set of pre-defined steps and organize them to form the plans for the tasks (Dalvi et al., 2021; Creswell and Shanahan, 2022; Qu et al., 2022; Hong et al., 2022; Creswell et al., 2022; Kazemi et al., 2023; Zhang et al., 2024). However, these studies are limited by the pre-defined steps set and customized for specific tasks, resulting in poor generalization ability. Although some studies have made initial attempts to conduct planning from back to front (Li et al., 2024; BAAI, 2023; Wang et al., 2023a), they require a significant amount of trial and error before acquiring sufficient knowledge for planning retrieval (e.g., Optimus-1 (Li et al., 2024) accumulates historical experience through iterative interaction with the environment to support plan retrieval from the knowledge graph).

To alleviate the above defects, we design a backward reasoning based agent that is free from the impact of the big perception gap between the agent's initial state and task goal, and enable it to handle various tasks by decomposing the task goals into natural language steps and sub-goals with LLM. Moreover, our agent supports direct static planning for any complex task by recursive goal decomposition, achieving good results even without interacting with the environment.

2.2 Agent for Minecraft

Minecraft provides a simulation environment with a high degree of freedom and a wide range of tasks that can be decomposed and executed, making it an ideal resource for exploring the planning ability of LLM-based agents. Some studies rely on extensive interaction with the environment to gradually piece together the proper plans (Wang et al., 2023a; Zhu et al., 2023; Wang et al., 2023d; Liu et al., 2024). Other studies rely on a pre-designed capabilities library to generate the plans using humanwritten rules (Yuan et al., 2023; Wang et al., 2023c; Zhao et al., 2023). However, these studies require a significant amount of effort to trial and error and only support limited tasks with specific customization. Thus the trained models have high coupling with specific tasks and cannot generalize to other tasks. To address these defects, we first enable our agent to support static planning with LLMs that

do not require extra customization for particular tasks, and then equip it with the ability to improve the efficiency of planning by utilizing the agent's interaction with the environment.

3 Method

Our agent **BAR** consists of three key modules: (a) recursive goal decomposition module; (b) state consistency maintaining module; (c) stage memory module. The first two modules are designed to provide robust and consistent planning, and the third module is designed to further improve planning efficiency by interacting with the environment.

3.1 Recursive Goal Decomposition

To tackle the challenge brought by the big perception gap between the agent's initial state and task goal, we propose to conduct the planning starting from the terminal state (or goal state). Moreover, since decomposing one task into steps in a single iteration is extremely challenging, we propose to decompose each task through multiple iterations. Drawing inspiration from dynamic programming (Bellman, 1954), we introduce sub-goals as the transition between the task goal and the final generated plan to alleviate this planning difficulty in a recursive manner.

In the dynamic programming algorithm (Bellman, 1954), the solution to the current problem in each iteration is composed of two parts --- (1) the immediately obtainable result and (2) the remaining results that need further solving in subsequent iterations. Inspired by this, we design a recursive goal decomposition module to split each task goal into one step and some sub-goals that need further decomposition. As shown in Figure 3 (a), given the task goal "Obtain 1 🄊 (diamond pickaxe)", in the first iteration we decompose the goal into two parts: (1) a step needs to be executed to achieve the goal ("Craft 1 \nearrow ") and (2) sub-goals that are necessary to support the execution of the decomposed step ("Collect 3 ", "Obtain 2 /", and "Obtain 1 **()**").

After finishing one iteration of such a decomposition, the decomposed step is pushed into a "Step Stack" that stores all the steps to form the final plan, and the decomposed sub-goals are sequentially queued into a "Goal Queue" that stores the remaining goals that need to be decomposed in sub-



Figure 3: Our proposed backward reasoning based agent. The agent consists of three key modules. Module (a) is used to generate robust plans through recursive goals decomposition. Module (b) is used to eliminate the state conflicts in different steps by integrating forward and backward reasoning. Module (c) is used to further improve the efficiency of planning by interacting with the environment. Best viewed on screen.

sequent iterations. In the next iteration, we fetch a new goal from the head of the Goal Queue and repeat the above decomposition process to the new goal. We continuously perform such decompositions until the Goal Queue is empty. Finally, we sequentially pop the steps from the Step Stack to obtain the final plan for the given task. Since there may be multiple steps of the same type, we fuse them into a single step by adding up the numbers of the items in these steps, for example, fuse "Obtain $3 \checkmark$ " and "Obtain $2 \checkmark$ " into "Obtain $5 \checkmark$ ". The pseudocode of the recursive goal decomposition module is shown in Algorithm 1.

Through this module, we avoid the big perception gap between the agent's initial state and task goal by decomposing the task goal starting from the terminal state. And we reduce the difficulty of planning through multiple iterations of decomposition. Moreover, the decomposed steps and sub-goals exist in the form of natural language, without relying on external predefined steps set or customization for specific tasks. This enables our agent to handle the planning of arbitrary tasks.

3.2 State Consistency Maintenance

Although recursive goal decomposition is not hindered by the big perception gap between the agent's initial state and task goal, the state conflicts in dif-

Algorithm 1 Recursive Goal Decomposition

Input: Task goal G

- 1: Initialization: goal_queue GQ = [G], step_stack SS = [], plan P = []
- 2: repeat
- 3: top_goal $\leftarrow GQ.get_top_goal()$
- 4: decomposed_step, sub_goals ← Decompose(top_goal)
- 5: *SS*.push(decomposed_step)
- 6: **for** *i* in range (len(sub_goals)) **do**
- 7: $GQ.push(sub_goals[i])$
- 8: end for
- 9: **until** GQ is empty
- 10: repeat
- 11: top_step $\leftarrow SS.pop()$
- 12: *P*.append(top_step)
- 13: **until** SS is empty

Output: P

ferent steps during planning makes it not yet a perfect solution. In fact, by forward reasoning, the agent can clearly understand its physical states before and after executing each step. Inspired by this, we propose to maintain state consistency during planning by integrating forward and backward reasoning.

To achieve this, after generating the initial plan by module (a), we choose pairs of start and end anchor steps from the initial plan, between which an incorrect plan may be generated due to the state conflicts between different steps. Then we adopt forward reasoning based planning to generate the partial plan between the chosen start and end anchor steps. Finally, we integrate the initial plan and the complementary partial plan to correct possible mistakes in the initial plan. We design two methods to choose the proper start and end anchor steps: (I) step scoring: we rate all the steps in the initial plan on a scale of 1 to 10, with 1 indicating the next few steps starting from this step are likely to be wrong and 10 indicating the next few steps starting from this step are completely accurate. Then we choose the steps with a rating below the threshold t as the start anchor steps and choose the k-th step after each start anchor step as the end anchor steps. (II) sliding window: we randomly select the start and end anchor steps with an interval of k steps from the initial plan.

As shown in Figure 3 (b), given the initial plan of task "Obtain 1 \nearrow ", by the step scoring method, we first rate all the steps and find that the rating of step "Craft 1 (crafting table)" is below the threshold t. So we choose it as the start anchor step and choose the k-th step after it ("Craft 1 (furnace)") as the end anchor step. Then we adopt forward reasoning based planning to generate a partial plan between them. After integrating the initial plan and the generated partial plan, the agent realizes that the initial plan ignores one important step: before mining (stone), it should reach below ground by digging down where exist. Therefore, the agent adds the missing step and obtains a refined plan.

Through this module, we force the agent to understand its physical states before and after executing each step and eliminate state conflicts that might occur between different steps during planning.

3.3 Stage Memory

In addition to utilizing the reasoning ability of LLMs, taking advantage of the interaction with the environment can further improve the planning efficiency of our agent. Current studies utilize the execution records of generated plans to support the planning of other more complex tasks (Zhu et al., 2023; Wang et al., 2023d; Lin et al., 2023). However, these studies record the execution results of entire plans that are not applicable to new tasks

with different task goals, resulting in inefficient memory utilization.

We solve this problem by designing a stage memory module. As our agent decomposes the task goal in multiple iterations, the decomposition result of each iteration can be used to guide subsequent goal decomposition for new tasks. To pick out accurate decomposition results, we execute the plans generated by modules (a) and (b) in the environment. The higher the execution success rate, the more accurate the decomposition result of the corresponding plan is. After executing the generated plans, we record each plan's execution success rate and corresponding goal decomposition results during the planning, called stage memory. When planning for new tasks, we retrieve the decomposition results of the same goals from stage memory with a high execution success rate and assist in decomposing the goals in new tasks.

As shown in Figure 3 (c), after executing generated plans of task "Obtain 1 f (fence)", "Collect 1 \checkmark (gold ingot)" and so on, we record each plan's goal decomposition results and execution success rate into the stage memory. When planning for a new task "Collect 1 \blacklozenge (redstone)", we can retrieve the decomposition result of the same goal "Obtain 1 \nearrow (iron pickaxe)" from stage memory and assist in decomposing the goal for this new task. Therefore, our agent can continuously improve the efficiency of goal decomposition with the help of stage memory. See appendix C for more details of our agent.



Figure 4: Collected tasks on the technology tree in Minecraft. As the depth increases, the length of the plans for the tasks gradually increases.

4 Experiments

To evaluate the effectiveness of our proposed agent **BAR**, we conduct extensive experiments to compare the planning performance of **BAR** with current SOTA baselines, on both static and dynamic planning settings.

4.1 Experimental Setup

Static Planning Setting: To compare the performance of different methods, we first conduct static planning experiments with textual metrics. We constructed a dataset consisting of various tasks on the inherent technology tree in Minecraft. In total, we collect 53 tasks across 5 distinct groups, ranging from stone group to gold group with varying plan lengths, as shown in Figure 4. We hired three experts with extensive experience in Minecraft to annotate the unique decomposed steps for each task as the ground truth, and all annotations were required to have been agreed by all three experts. We only compare the plans generated by modules (a) and (b) of our agent with other methods for static planning setting, without interacting with the environment. For choosing anchor steps in section 3.2, we set t to 5 and k to 3, and discuss their performances in section 4.4. Each task was performed 10 times by each model and the average results were recorded. The temperatures of all the LLMs used are set to 0 for reproducibility. Finally, we report the mean performance for each group's tasks.

Dynamic Planning Setting: To further evaluate the potential of our agent utilizing the interaction with the environment, we conduct dynamic planning experiments. We first execute the plans generated by modules (a) and (b) in section 3 10 times and record their average execution success rates into stage memory. Then we conduct planning again for all tasks with the help of the decomposition results corresponding to the plans with success rate ≥ 0.3 in stage memory. We report the average results of each group's tasks. See appendix B for more detailed experimental settings.

4.2 Static Planning

Baselines: (1) **Chain-of-thought** (Wei et al., 2022): utilize forward reasoning based thought chains to generate the plans step by step. (2) **Reverse Chain** (Zhang et al., 2024): utilize backward reasoning based thought chains to generate

the plans. (3) Self-Refine (Madaan et al., 2024): utilize feedback from the model itself to refine the generated plans. (4) Tree-of-thought (Yao et al., 2024): enable the exploration of multiple forward reasoning paths and choose the best one. (5) **DEPS** (Wang et al., 2023c): enhance the accuracy of planning by integrating the description and explanation of the plans given by LLM. (6) Plan-and-Solve (Wang et al., 2023b): utilize incontext learning to rehearse generated plans and improve the logical correctness of the plans. (7) Openai-o1 (OpenAI, 2024): perform a lot of reasoning and reflection before outputting the final results to improve the accuracy of generated plans. See appendix D for implementation details of the baselines.

Metrics: We design three textual metrics based on overlap computation to evaluate the quality of generated plans: (1) **Accuracy**: Evaluate whether each step in the generated plan matches exactly with the corresponding step in the ground truth at the same index. (2) **F1-Score**: Take into account both the precision and recall of generated plans. A step is considered correctly generated only if the pair it forms with its previous step can match the ground truth. (3) **Edit Distance**: Evaluate the minimum number of editing operations needed to transform the generated plan into the ground truth. Editing operations involve replacing, inserting, or deleting a single step in the generated plan. See appendix B.2 for detailed demonstration.

Results: As shown in Table 1, our agent achieves SOTA performance on all 5 groups under all three metrics. This demonstrates that our agent can provide robust and consistent planning for various tasks. In addition, we have the following findings:

(1) In addition to the powerful closed-source model GPT4 (Achiam et al., 2023), with the opensource Llama-3 model (Dubey et al., 2024) as backbone, our agent can still significantly outperform other baselines that also adopt Llama-3 as backbone model. This shows the robustness of our agent for planning with both closed-source and open-source LLMs.

(2) Chain-of-thought fails to achieve good performance due to the inherent challenge brought by the big perception gap between the agent's initial state and task goal. This makes it difficult for an LLM to infer correct steps, and the situation becomes worse on tasks with longer plans.

Method	stone			iron			diamond			redstone			gold		
	Acc ↑	F1 ↑	$\text{ED}\downarrow$	Acc \uparrow	F1 ↑	$\text{ED}\downarrow$									
Chain-of-thought															
- Llama-3-8B	68.97	68.13	2.91	74.73	74.95	3.29	55.21	51.47	8.29	43.69	48.45	8.57	41.28	52.75	9.00
- GPT4	75.82	75.00	2.09	77.26	78.12	2.86	68.47	71.94	4.29	47.66	54.88	7.29	52.54	73.64	5.86
Reverse Chain															
- Llama-3-8B	22.50	44.18	5.36	53.28	65.01	4.67	29.47	35.62	9.57	13.64	20.84	12.43	23.81	21.23	12.86
- GPT4	73.33	74.03	2.18	80.87	81.77	2.43	86.49	86.49	2.14	77.88	77.88	3.57	52.17	74.56	5.43
Self-Refine															
- Llama-3-8B	65.85	63.70	3.18	69.23	69.58	4.05	40.62	48.40	10.29	40.57	47.70	8.57	37.27	48.92	9.57
- GPT4	75.82	75.00	2.09	74.01	77.80	3.10	74.77	79.18	3.29	56.88	66.64	5.71	52.54	74.09	5.71
Tree-of-thought	46.91	43.92	4.27	33.96	44.46	6.81	38.64	28.43	10.00	20.43	22.23	12.14	35.71	30.95	10.14
DEPS	49.45	46.20	3.73	61.01	50.35	5.24	49.55	41.32	7.86	26.44	25.95	11.00	48.65	39.42	8.14
Plan-and-Solve	45.45	58.29	3.45	76.17	80.29	2.71	66.67	77.02	4.14	42.06	50.67	8.43	63.00	59.93	6.43
Openai-o1	74.73	74.86	2.09	77.98	82.34	2.67	63.96	79.55	3.57	57.52	76.44	3.86	75.42	83.97	3.14
Ours															
- Llama-3-8B	78.02	78.02	1.82	74.44	76.00	3.33	72.97	79.67	3.29	59.22	59.10	7.00	60.71	66.38	5.86
- GPT4	87.91	87.91	1.00	82.67	85.02	2.19	87.39	88.23	1.86	80.53	82.22	2.86	85.59	87.34	2.14

Table 1: Results of the static planning experiments. Best results are highlighted in **bold**, and the second best results are highlighted with <u>underline</u>. For the adopted metrics, "Acc" refers to Accuracy, "F1" refers to F1-Score, and "ED" refers to Edit Distance. Larger Accuracy and F1-Score and smaller Edit Distance represent better performance.

(3) Backward reasoning shows its superiority by comparing the performance of Chain-of-thought and Reverse Chain. Contrary to Chain-of-thought, Reverse Chain infers the steps to be executed starting from the terminal state, free from the big perception gap between the agent's initial state and task goal. However, simply changing the reasoning paradigm cannot solve the defects brought by the state conflicts in different steps during planning, resulting in a decline in performance.

(4) Self-Refine and DEPS can improve the plans by obtaining feedback from the LLM. However, as the initial plans are generated based on forward reasoning, they face the same challenge as Chainof-thought. Moreover, severe error propagation during multiple turns of refinement may introduce irrelevant noise to the generated plans, and the unstable feedback from the LLM can mislead the agent's reasoning to the wrong direction.

(5) Although Tree-of-thought can explore multiple paths during planning, evaluating the quality of incomplete plans by LLM is unreliable. Moreover, due to the interference of a large number of candidate steps, the generated plans are mixed with unnecessary noisy steps, resulting in poor performance.

(6) Plan-and-Solve and Openai-o1 can improve the plans by previewing the decomposed partial plans and repeatedly considering more possible solutions. However, even with more exploration of different paths and reflection on generated plans, it is still difficult for an LLM to infer correct steps that are far from the task goal.

4.3 Dynamic Planning

To further evaluate the potential of our agent utilizing the interaction with the environment, with the help of our designed stage memory described in section 3.3, we compare our agent with other baselines that also utilize the interaction with the environment for planning.

Baselines: (1) **React** (Yao et al., 2023): combine reasoning and acting to conduct planning. (2) **Jarvis-1** (Wang et al., 2023d): conduct planning in Minecraft with both pre-trained knowledge in LLMs and the game experiences given by the environment. (3) **Ours-Static**: our agent with only modules (a) and (b) in static planning setting, GPT4 as the backbone model. (4) **React + Ours-Static**: execute the generated plans by **Ours-Static**, but under the reasoning and acting framework of **React**.

Metrics: In addition to the textual metrics used in the static planning experiments, we further record the success rates and computation time for executing the generated plans in the environment (execution metrics). For a fair comparison, we execute all evaluated tasks 10 times for and report the average execution success rates and computation time for tasks in each group.

Results: As shown in Table 2 and Figure 5, with the support of the stage memory module, our agent beats the SOTA performance by utilizing the interaction with the environment, both in terms of textual metrics and execution metrics. Additionally, we have the following observations:

(1) Based on our designed stage memory module,

Method	stone			iron			diamond			redstone			gold		
	$Acc\uparrow$	F1 ↑	$ED\downarrow$	Acc \uparrow	F1 ↑	$\text{ED}\downarrow$	Acc \uparrow	F1 ↑	$\text{ED}\downarrow$	Acc \uparrow	F1 ↑	$\text{ED}\downarrow$	Acc \uparrow	F1 ↑	$\mathrm{ED}\downarrow$
React	70.33	72.53	2.18	89.17	89.17	1.43	71.17	74.77	4.57	57.01	57.49	7.00	57.63	79.75	4.43
Jarvis-1	75.82	75.00	2.09	73.61	76.56	3.38	81.98	81.98	2.86	70.09	71.70	5.00	52.54	74.09	5.71
Ours															
- static	87.91	87.91	1.00	82.67	85.02	2.19	87.39	88.23	1.86	80.53	82.22	2.86	85.59	87.34	2.14
- dynamic	92.31	92.31	0.64	96.58	96.92	0.45	88.29	89.23	1.71	86.73	86.73	2.14	86.44	94.93	0.86

Table 2: Results of the dynamic planning experiments by textual metrics. Best results are highlighted in **bold**. For the adopted metrics, "Acc" refers to Accuracy, "F1" refers to F1-Score, and "ED" refers to Edit Distance. "Ours-static" refers to our agent with only modules (a) and (b) in static planning setting, GPT4 as the backbone model. "Ours-dynamic" refers to our complete agent, including the stage memory module, GPT4 as the backbone model.

our agent can obtain the feedback of generated plans by interacting with the environment, thereby improving the efficiency of subsequent goals decomposition. With just one round of interaction, the quality of the generated plans can be significantly improved. (2) React and Jarvis-1 can achieve good performance by interacting with the environment. However, they cannot solve the big perception gap between the agent's initial state and task goal during planning. Even if the feedback from the environment can indicate the execution results of generated plans, no effective information is given to refine the incorrect plans. (3) Taking the plans generated by our agent as the reference, React can significantly improve the execution success rates of plans. This indicates that high-quality plans are key in helping agents successfully complete tasks in the environment. Therefore, relying solely on environmental feedback to correct the generated plans is not the most effective solution. A better approach is to focus on improving the quality of the generated initial plans by LLM. (4) Environmental feedback can indeed help improve the quality and execution success rate of plans. But compared with React and Jarvis-1, our agent can more effectively utilize the interaction with the environment through stage memory, as we can efficiently take advantage of every piece of decomposition results in all tasks.

4.4 Ablation Study

The effectiveness of our proposed stage memory module has been verified in section 4.3. Now we conduct an ablation study to verify the effectiveness of the state consistency maintaining module and compare the two methods to choose anchor steps described in section 3.2. As shown in Figure 6 and 7, our proposed state consistency maintaining module can significantly improve the quality of generated plans. Whether using GPT4 or Llama-3 as the backbone model, it can effectively eliminate



Figure 5: Results of the dynamic planning experiments by execution metrics. The bars represent the average execution success rates of tasks in different groups, and the lines represent the average time spent to execute the takes.

the state conflicts in different steps by integrating forward and backward reasoning. The open-source Llama-3 model can benefit more from it because the initial plans' quality of Llama-3 is lower than that of GPT4. Additionally, by comparing the two methods to select anchor steps, we find that step scoring can bring greater performance improvement. This is because the sliding window method chooses too many anchor step pairs and introduces unnecessary noise into the generated plans.



Figure 6: Results of the ablation study for state consistency maintaining module, GPT4 as the backbone model.

5 Conclusion

In this paper, we study the planning problem for complex tasks in the Minecraft environment. We



Figure 7: Results of the ablation study for state consistency maintaining module, Llama-3-8B as the backbone model.

argue that, compared with forward reasoning, backward reasoning has the advantage to avoid the perception gap between the agent's initial state and the task goal. We further point out a key weakness of backward reasoning is that it may fail to complete the task because of the state conflicts during planning. To this end, we propose **BAR** with: (1) recursive goal decomposition to make robust backward reasoning, (2) a state consistency maintaining module to address state conflicts, (3) and a stage memory module promote more efficient planning. Experimental results show the superiority of our agent in both static and dynamic planning, surpassing other baselines by a large margin.

6 Ethical Considerations

The dataset we constructed does not contain any personal, private or sensitive information and is only used for research purposes. The participation of volunteer human annotators in the construction of the dataset does not constitute an ethical concern, as all annotators were informed of the nature of the task, participated voluntarily without coercion, and no sensitive personal data was collected or disclosed during the annotation process. Therefore, we believe that there is no ethical issue with our work.

7 Limitations

To verify the generalizability of our work, tests on many further tasks should be undertaken, but we believe that as a primary step towards Backward Reasoning based Planning with LLMs, Minecraft is rich enough to provide many scenarios and tasks in an open-ended environment that simulates the real world. As described in section 4.1, we have collected a large number of tasks with various difficulty levels in different scenarios to demonstrate that our method can handle a wide range of tasks. Moreover, as demonstrated in prior work, Minecraft is a sufficiently complex environment to undertake a rich variety of planning experiments. We leave exploration in other environments as future work.

Acknowledgement

This work was supported in part by the National Natural Science Foundation of China (No. 62206191, No. 62272330, and No. U24A20328); in part by the Science Fund for Creative Research Groups of Sichuan Province Natural Science Foundation (No. 2024NSFTD0035); in part by the Natural Science Foundation of Sichuan (No. 2024YFHZ0233).

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. GPT-4 technical report. arXiv preprint arXiv:2303.08774.
- P BAAI. 2023. Plan4mc: Skill reinforcement learning and planning for open-world minecraft tasks. *arXiv preprint arXiv:2303.16563*.
- Richard Bellman. 1954. The theory of dynamic programming. *Bulletin of the American Mathematical Society*, 60(6):503–515.
- Ziru Chen, Michael White, Ray Mooney, Ali Payani, Yu Su, and Huan Sun. 2024. When is tree search useful for LLM planning? it depends on the discriminator. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13659–13678, Bangkok, Thailand. Association for Computational Linguistics.
- Antonia Creswell and Murray Shanahan. 2022. Faithful reasoning using large language models. *arXiv preprint arXiv:2208.14271*.
- Antonia Creswell, Murray Shanahan, and Irina Higgins. 2022. Selection-inference: Exploiting large language models for interpretable logical reasoning. *arXiv* preprint arXiv:2205.09712.
- Bhavana Dalvi, Peter Jansen, Oyvind Tafjord, Zhengnan Xie, Hannah Smith, Leighanna Pipatanangkura, and Peter Clark. 2021. Explaining answers with entailment trees. In *Proceedings of the 2021 Conference* on Empirical Methods in Natural Language Processing, pages 7358–7370.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,

Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The Llama 3 herd of models. *arXiv* preprint arXiv:2407.21783.

- Ruixin Hong, Hongming Zhang, Xintong Yu, and Changshui Zhang. 2022. Metgen: A module-based entailment tree generation framework for answer explanation. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 1887– 1905.
- Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. 2022. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *International conference on machine learning*, pages 9118–9147. PMLR.
- Xu Huang, Weiwen Liu, Xiaolong Chen, Xingmei Wang, Hao Wang, Defu Lian, Yasheng Wang, Ruiming Tang, and Enhong Chen. 2024. Understanding the planning of LLM agents: A survey. *arXiv* preprint arXiv:2402.02716.
- Weisen Jiang, Han Shi, Longhui Yu, Zhengying Liu, Yu Zhang, Zhenguo Li, and James Kwok. 2024. Forward-backward reasoning in large language models for mathematical verification. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 6647–6661.
- Mehran Kazemi, Najoung Kim, Deepti Bhatia, Xin Xu, and Deepak Ramachandran. 2023. Lambada: Backward chaining for automated reasoning in natural language. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pages 6547–6568.
- Jing Yu Koh, Stephen McAleer, Daniel Fried, and Ruslan Salakhutdinov. 2024. Tree search for language model agents. *arXiv preprint arXiv:2407.01476*.
- Zaijing Li, Yuquan Xie, Rui Shao, Gongwei Chen, Dongmei Jiang, and Liqiang Nie. 2024. Optimus-1: Hybrid multimodal memory empowered agents excel in long-horizon tasks. *arXiv preprint arXiv*:2408.03615.
- Jiaju Lin, Haoran Zhao, Aochi Zhang, Yiting Wu, Huqiuyue Ping, and Qin Chen. 2023. Agentsims: An open-source sandbox for large language model evaluation. *arXiv preprint arXiv:2308.04026*.
- Shaoteng Liu, Haoqi Yuan, Minda Hu, Yanwei Li, Yukang Chen, Shu Liu, Zongqing Lu, and Jiaya Jia. 2024. Rl-GPT: Integrating reinforcement learning and code-as-policy. arXiv preprint arXiv:2402.19299.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. 2024. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36.
- Drew McDermott. 1991. Regression planning. International Journal of Intelligent Systems, 6(4):357–416.

- Allen Newell, John C Shaw, and Herbert A Simon. 1959. Report on a general problem solving program. In *IFIP congress*, volume 256, page 64. Pittsburgh, PA.
- OpenAI. 2024. Learning to reason with LLMs. https://openai.com/index/learning-to-reason-with-llms/.
- Hanhao Qu, Yu Cao, Jun Gao, Liang Ding, and Ruifeng Xu. 2022. Interpretable proof generation via iterative backward reasoning. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2968–2981.
- Soumya Sanyal, Harman Singh, and Xiang Ren. 2022. Fairr: Faithful and robust deductive reasoning over natural language. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1075–1093.
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. 2023a. Voyager: An open-ended embodied agent with large language models. *Transactions on Machine Learning Research*.
- Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. 2023b. Plan-and-solve prompting: Improving zeroshot chain-of-thought reasoning by large language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pages 2609–2634.
- Zihao Wang, Shaofei Cai, Guanzhou Chen, Anji Liu, Xiaojian Ma, Yitao Liang, and Team CraftJarvis. 2023c. Describe, explain, plan and select: interactive planning with large language models enables open-world multi-task agents. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, pages 34153–34189.
- Zihao Wang, Shaofei Cai, Anji Liu, Yonggang Jin, Jinbing Hou, Bowei Zhang, Haowei Lin, Zhaofeng He, Zilong Zheng, Yaodong Yang, et al. 2023d. Jarvis-1: Open-world multi-task agents with memoryaugmented multimodal language models. arXiv preprint arXiv:2311.05997.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824–24837.
- Tianci Xue, Ziqi Wang, Zhenhailong Wang, Chi Han, Pengfei Yu, and Heng Ji. 2023. RCoT: Detecting and rectifying factual inconsistency in reasoning by reversing chain-of-thought. *arXiv preprint arXiv:2305.11499*.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2024. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36.

- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*.
- Fangxu Yu, Lai Jiang, Haoqiang Kang, Shibo Hao, and Lianhui Qin. 2024a. Flow of reasoning: Efficient training of LLM policy with divergent thinking. *arXiv preprint arXiv:2406.05673*.
- Fei Yu, Hongbo Zhang, Prayag Tiwari, and Benyou Wang. 2024b. Natural language reasoning, a survey. *ACM Computing Surveys*.
- Haoqi Yuan, Chi Zhang, Hongcheng Wang, Feiyang Xie, Penglin Cai, Hao Dong, and Zongqing Lu. 2023. Skill reinforcement learning and planning for open-world long-horizon tasks. *arXiv preprint arXiv:2303.16563*.
- Yinger Zhang, Hui Cai, Xierui Song, Yicheng Chen, Rui Sun, and Jing Zheng. 2024. Reverse chain: A generic-rule for LLMs to master multi-API planning. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 302–325.
- Zhonghan Zhao, Wenhao Chai, Xuan Wang, Li Boyi, Shengyu Hao, Shidong Cao, Tian Ye, Jenq-Neng Hwang, and Gaoang Wang. 2023. See and think: Embodied agent in virtual environment. *arXiv preprint arXiv:2311.15209*.
- Enshen Zhou, Yiran Qin, Zhenfei Yin, Yuzhou Huang, Ruimao Zhang, Lu Sheng, Yu Qiao, and Jing Shao. 2024. Minedreamer: Learning to follow instructions via chain-of-imagination for simulated-world control. *arXiv preprint arXiv:2403.12037*.
- Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Weijie Su, Chenyu Yang, Gao Huang, Bin Li, Lewei Lu, Xiaogang Wang, et al. 2023. Ghost in the Minecraft: Generally capable agents for open-world environments via large language models with textbased knowledge and memory. *arXiv preprint arXiv:2305.17144*.

A Introduction to the Simulation Environment Minecraft

Minecraft is a simulation environment with extremely high degrees of freedom and provides a large number of scenarios and tasks similar to the real world for agents to explore. In Minecraft, intelligent agents can do anything they want in different scenarios, such as collecting various materials and crafting useful tools. Therefore, exploring the behaviors of agents in Minecraft is of great significance for serving various applications in the real world. Minecraft covers diverse types of tasks that can be executed, with sufficiently complex tasks to challenge the agents' abilities. Moreover, the simulation environment is completely under our control, which facilitates us to conduct a large number of scientific experiments. Based on these advantages, we chose Minecraft as the platform to explore the automatic task completion and planning abilities of LLM based agents.

B Details of Experimental Settings

B.1 Details of Evaluated Tasks

To evaluate the planning ability of different models, we constructed a dataset consisting of various tasks on the technology tree in Minecraft.

Considering that the evaluation metrics are based on overlap computation to assess the accuracy of generated plans, we only selected tasks with unique plans when constructing the dataset. This can make sure that the ground truth for each task is unique. For a few flexible tasks, we first standardize the plans generated by the model according to our predefined unified standards and then compare them with the ground truth to ensure fair comparison (e.g., always ensure that the step "craft 1 crafting_table" is before the step "craft 1 stick").

Details of the chosen tasks for evaluation are listed in Table 3.

Table 3: Details of Evaluated Tasks

Group	Task Num	Avg Plan Length	Example
Stone	11	8.27	obtain 1 🛢
Iron	21	13.27	obtain 1 🖉
Diamond	7	15.86	obtain 1 🗶
Redstone	7	16.14	obtain 1 🜑
Gold	7	16.86	obtain 1 🏓

B.2 Details of Adopted Textual Metrics

We design three textual metrics based on overlap computation to evaluate the quality of generated plans. The computation demonstration of the three metrics is shown in Figure 8.



Figure 8: Computation demonstration of the three adopted textual metrics.

B.3 Details of the Simulation Environment

B.3.1 High-level Steps Encapsulation

As our work focuses on task planning instead of generating fine-grained machine instructions, based on the work of Jarvis-1 (Wang et al., 2023d), we utilize our game experience and knowledge of Minecraft to encapsulate complex operations into high-level steps that align with the steps in generated plans by our agent and baseline models. The action space in Minecraft includes keyboard and mouse operations that are usually performed by humans. For example, using the keyboard to control the agent's movement, jumping, opening or closing the inventory, and selecting items from the hotbar. Mouse movements are used to adjust the agent's view, including horizontal rotation and pitch. Mouse clicks are used for attacking mobs, breaking blocks, and interacting with crafting_table and furnace. When using crafting_table and furnace, precise item selection and manipulation are needed with fine-grained mouse movement control, which cannot be reliably accomplished by current low-level controller models.

To address this problem and make sure the generated plans can be successfully executed in the environment, we propose to encapsulate the finegrained keyboard and mouse operations into three types of frequently-used steps: craft, smelt, and equip. Heuristic rules are applied to the three types of steps for encapsulation. For another two types of important steps mine and dig down, we use the SOTA low-level controller model Mine-Dreamer(Zhou et al., 2024) to directly execute the natural language instructions. Details of the action space in Minecraft and the encapsulation of the above mentioned steps are listed in Table 4 and Table 5.

B.3.2 Environment Setting

In dynamic planning experiments, to make sure the comparison between different methods is fair and the agent can successfully execute the generated plans, we need to control the environment to remove the influence of irrelevant factors in the environment. For example, for the task "Craft 1 \nearrow (diamond pickaxe)", when executing the step "Mine 3 (diamond)", the agent may not be able to mine enough diamonds even if the agent uses the right tool to mine in below the ground because the storage of diamonds below the ground is too scarce.

To solve this problem, we conduct two types of initialization in the environment:

Global Environment Initialization. The rules used to set global setting are listed as follows:

- We set the agent to be born in a fixed position in a fixed scene with an empty inventory each time. Therefore, the agent executes the plan starting from a fixed initial state every time.
- We set the agent to have a bright field of view even in the dark. Therefore, the execution of steps will not be affected by the darkness when the agent is below the ground.
- We set the environment to peaceful mode so that

no other creatures (such as zombies) can interfere with the agent to execute the plans.

• We set that the environment is always in daylight. So the agent can better execute the plans on above the ground with adequate lighting.

Ore Distribution Initialization. To ensure that there is sufficient ore below the ground for the agent to mine, we place some ore below the ground when initializing the environment, so that the agent is more likely to mine the required ore as long as the plans are correct. The command used to place the ore is listed as follows:

/fill <from_x> <from_y> <from_z> <to_x> <to_y> <to_z> minecraft:ore

Specifically, we implemented an ore distribution algorithm to place the ore. This algorithm generates multiple concentric square layers of ore around the agent's born position in below the ground. In this way, there will be enough ore below the ground for the agent to mine.

C Details of Our Agent

We list the prompts used in our designed agent as follows.

To conduct the goal decomposition by LLM, we prompt the LLM to decompose each goal into step and sub-goals. For a fair comparison, keeping consistent with other baselines, we only provide the recipe information of the target item required by the task goal, excluding the intermediate sub-goals. Prompt used for goal decomposition in section 3.1 is as follows:

System:

You are a helpful assistant in Minecraft. I will give you a goal to achieve in Minecraft, and you need to decompose the goal into a single step and a list of sub goals to achieve. Output the reasoning thought and the decomposed result. You can follow the history dialogue to make the decomposition.

User:

Goal: collect 3 stone. Thought:

Description
Move forward. Press "W" in keyboard.
Move backward. Press "S" in keyboard.
Move left. Press "A" in keyboard.
Move right. Press "D" in keyboard.
Jump up. Press "Space" in keyboard.
Open or close the inventory GUI. Press "E" in keyboard.
Attack items. Click left button in mouse.
Use items in front of the agent. Click right button in mouse.
Move agent's view in the horizontal direction.
Move agent's view in the vertical direction.
Select the item in the hotbar to equip it to the agent's hand. Press 1-9 in keyboard.

Table 4: Details of the action space in Minecraft.

Table 5: Details of the encapsulation of used steps in our experiments.

Step	Description
Craft	Take the crafting table from the inventory. Place it on the ground. Open the crafting table,
	select the required materials based on the recipe to craft the target item. Take the crafted
	target item. After that, close the crafting table, destroy the crafting table, and return it to the
	inventory.
Smelt	Take the furnace from the inventory. Place it on the ground. Open the furnace, select the
	required materials based on the recipe to smelt the target item. Take the smelted target item.
	After that, close the furnace, destroy the furnace, and return it to the inventory.
Equip	Locate the cell of the target item in the inventory or hotbar and equip it to the agent's hand.
Mine	Directly execute the natural language instructions using the SOTA low-level controller model
	MineDreamer (Zhou et al., 2024).
Dig down	Directly execute the natural language instructions using the SOTA low-level controller model
	MineDreamer (Zhou et al., 2024).

To collect 3 stone, the last step is to mine 3 stone with wooden_pickaxe, as mining the stone requires at least wooden_pickaxe. And the previous sub goals are to obtain 1 wooden_pickaxe and dig down with wooden_pickaxe, because stone only appears at the below ground level. Based on these analysis, the decomposed result is as follows:

Decomposed Step:

Mine 3 stone with wooden_pickaxe Decomposed Sub Goals:

1. Obtain 1 wooden_pickaxe

2. Dig down with wooden_pickaxe

=======

Goal: {new goal} Thought:

Assistant:

Prompt used to rate for all the steps in the initial plan in section 3.2 is as follows:

System:

You are a helpful assistant in Minecraft. I will give you a goal to achieve in Minecraft and the generated initial plan to achieve this goal. You need to rate for all the steps in the initial plan to evaluate the correctness of the steps on a scale of 1 to 10, where 1 indicating the next few steps starting from this step is likely to be wrong and 10 indicating the next few steps starting from this step is completely accurate.

User:

Goal: collect 3 stone.Initial Plan:1. Mine 3 log with barehand2. Craft 9 planks3. Craft 2 stick

4. Craft 1 crafting_table 5. Craft 1 wooden pickaxe 6. Mine 3 stone with wooden_pickaxe Thought: As the stone only exists in below the ground, when executing the step "Mine 3 stone with wooden_pickaxe", the agent should be below the ground. However, before this step the agent is gathering materials and crafting items above the ground. So the latter half of the steps in the initial plan may not be executed successfully. Rating: 1. Mine 3 log with barehand - 10 2. Craft 9 planks - 8 3. Craft 2 stick - 5

- 4. Craft 1 crafting table 3
- 5. Craft 1 wooden_pickaxe 3
- 6. Mine 3 stone with wooden_pickaxe 5

Goal: {goal of new task} Initial Plan: {generated initial plan} Thought:

Assistant:

Prompt used to complement the partial plan between the given start and end anchor steps in section 3.2 is as follows:

System:

You are a helpful assistant in Minecraft. I will give you a task goal to achieve in Minecraft and a pair of start and end anchor steps chosen from the corresponding plan to achieve the task goal. You need to complement the partial plan between the start and end anchor steps to help achieve the task goal.

User:

Goal: collect 3 stone.

Start Anchor Step: Craft 1 crafting_table End Anchor Step: Mine 3 stone with wooden_pickaxe

Thought:

To achieve the task goal "collect 3 stone", starting from the step "Craft 1 crafting_table", next I need to craft 1 wooden pickaxe to mine stone. But before mining stone, I need to reach below ground first as stone only exists underground. Finally I can mine 3 stone with wooden_pickaxe. Based on these analysis, the partial plan should be as follows: Partial Plan: Craft 1 crafting_table Craft 1 wooden_pickaxe Dig down with wooden_pickaxe Mine 3 stone with wooden_pickaxe _____

Goal: {goal of new task} Start Anchor Step: {given start anchor step} End Anchor Step: {given end anchor step}

Thought:

Assistant:

Prompt used to integrate the initial plan and the complementary partial plan is as follows:

System:

You are a helpful assistant in Minecraft. I will give you an initial plan to achieve a goal in Minecraft and a complementary partial plan. You need to compare the two plans and correct possible mistakes in the initial plan.

User:

Goal: collect 3 stone.

Initial Plan:

- 1. Mine 3 log with barehand
- 2. Craft 9 planks
- 3. Craft 2 stick
- 4. Craft 1 crafting_table
- 5. Craft 1 wooden_pickaxe
- 6. Mine 3 stone with wooden_pickaxe

Complementary Partial Plan:

- 4. Craft 1 crafting_table
- 5. Craft 1 wooden_pickaxe
- 6. Dig down with wooden_pickaxe
- 7. Mine 3 stone with wooden_pickaxe

Thought:

The initial plan ignores one important step "Dig down with wooden_pickaxe" that is presented in the complementary partial plan. As the stone is in below ground, the agent should reach below ground first and then mine the stone.

Corrected Plan:

- 1. Mine 3 log with barehand
- 2. Craft 9 planks
- 3. Craft 2 stick
- 4. Craft 1 crafting_table
- 5. Craft 1 wooden_pickaxe
- 6. Dig down with wooden_pickaxe
- 7. Mine 3 stone with wooden_pickaxe

Goal: {goal of new task}

Initial Plan: {generated initial plan}

Complementary Partial Plan: {complementary partial plan}

Thought:

Assistant:

D Details of Baseline Models

In this section, we will provide the details of the baselines used in our experiments.

D.1 Baselines in Static Planning

Baselines in static planning do not interact with the environment. The plans are generated by prompting LLMs with specific modeling methods in the original studies. We adopt "Meta-Llama-3-8B-Instruct" ² and "OpenAI-gpt-4" as the backbone

²https://huggingface.co/meta-llama/ Meta-Llama-3-8B-Instruct model for "Chain-of-thought", "Reverse Chain", and "Self-Refine".

Chain-of-Thought (Wei et al., 2022). Chain-of-Thought utilizes forward reasoning based thought chains to generate the plans step by step. Prompt used for Chain-of-Thought is as follows:

System:

You are a helpful assistant in Minecraft. I will give you a goal to achieve in Minecraft, and you need to decompose the goal into a sequence of steps to help achieve the goal. Output the reasoning thought and the decomposed result. You can follow the history dialogue to make the decomposition.

User:

======= Goal: collect 3 stone.

Recipe: stone can be mined below the ground with at least wooden_pickaxe.

Thought:

Now let's think step by step. I should first mine 3 log with barehand. Then use 3 log to craft 9 planks. Then use 2 planks to craft 2 stick. Then use 4 planks to craft 1 crafting_table. Then use 3 planks and 2 stick to craft 1 wooden_pickaxe on 1 crafting_table. Then dig down with wooden_pickaxe. Then mine 3 stone with wooden_pickaxe. Based on the above reasoning, the plan to achieve the goal is: Plan:

- 1. Mine 3 log with barehand
- 2. Craft 9 planks
- 3. Craft 2 stick
- 4. Craft 1 crafting_table
- 5. Craft 1 wooden_pickaxe
- 6. Dig down with wooden_pickaxe
- 7. Mine 3 stone with wooden_pickaxe

Goal: {goal of new task}. Recipe: {recipe of new task}.

Thought:

Assistant:

Reverse Chain (Zhang et al., 2024). In contrast to Chain-of-thought, Reverse Chain utilizes backward reasoning based thought chains to generate the plans. Prompt used in Reverse Chain is as follows:

System:

You are a helpful assistant in Minecraft. I will give you a goal to achieve in Minecraft, and you need to decompose the goal into a sequence of steps. Output the reasoning thought and the decomposed result. You can follow the history dialogue to make the decomposition.

User:

Goal: collect 3 stone.

Recipe: stone can be mined below the ground with at least wooden_pickaxe.

Thought:

Now let's think step by step by backward The final goal is to collect reasoning. 3 stone. To collect 3 stone, I need to mine 3 stone with wooden_pickaxe. stone usually exists underground, so I need to dig down with wooden pickaxe. То obtain 1 wooden_pickaxe, I need to craft 1 wooden pickaxe using 3 planks and 2 stick on 1 crafting table. To obtain 1 crafting_table, I need to craft 1 crafting_table using 4 planks. To obtain 2 stick, I need craft 2 stick using 2 planks on 1 crafting table. So I totally need 9 planks. To obtain 9 planks, I need to craft 9 planks using 3 log, so I need to mine 3 log first. Mining 3 log can be done with barehand, so I start from here. Based on the above reasoning, the plan to complete this task is: Plan:

1. Mine 3 log with barehand

- 2. Craft 9 planks
- 3. Craft 2 stick
- 4. Craft 1 crafting_table
- 5. Craft 1 wooden_pickaxe
- 6. Dig down with wooden_pickaxe
- 7. Mine 3 stone with wooden_pickaxe

Goal: {goal of new task}. Recipe: {recipe of new task}.

Thought:

Assistant:

Self-Refine (Madaan et al., 2024). Self-Refine first generates an initial plan using the same prompt as Chain-of-Thought. The generated initial plan is then sent to the same LLM to obtain feedback. Finally, the feedback and the generated initial plan are sent to the same LLM to produce a refined plan. This process is repeated until the LLM's feedback indicates that the generated plan is totally correct.

Prompt used to generate the feedback for the already generated plan is as follows:

System:

You are a helpful assistant in Minecraft. I will give you a goal and the corresponding plan to achieve it in Minecraft, and you need to check the correctness of the plan and output the feedback. You need to output "Correct" or "Incorrect" first. If outputting "Incorrect", then output the detailed analysis and list the mistakes.

User:

Goal: collect 3 stone.

Recipe: stone can be mined below the ground with at least wooden_pickaxe.

Plan:

- 1. Mine 3 log with barehand
- 2. Craft 9 planks
- 3. Craft 2 stick
- 4. Craft 1 wooden_pickaxe
- 5. Dig down with wooden_pickaxe
- 6.Mine 3 log with wooden_pickaxe
- 7. Mine 3 stone with wooden_pickaxe

What's wrong with the plan? Please give me the feedback and list the mistakes if the plan is wrong.

Incorrect. Step 4 "Crafting 1 wooden_pickaxe" must be done on a crafting_table. But the plan does not include "Craft 1 crafting_table" before Step 4, so the first mistake is "ignoring obtaining one crafting_table before using it to craft other tools". Step 6 "Mine 3 log with wooden_pickaxe" is executed after Step 5 "Dig down with wooden_pickaxe". But log only exists on above the ground. So the second mistake is "collecting materials on above the ground when the agent is already in below the ground".

Mistake:

1. Ignoring obtaining one crafting_table before using it to craft other tools.

2. Collecting materials on above the ground when the agent is already in below the ground.

Goal: {goal of new task} Recipe: {recipe of new task}

Plan: {generated initial plan of new task}

What's wrong with the plan? Please give me the feedback and list the mistakes if the plan is wrong.

Assistant:

Prompt used to refine the initial plan based on the feedback given by LLM:

System:

You are a helpful assistant in Minecraft. I will give you a goal to achieve in Minecraft, the initial plan of this goal and the feedback about the initial plan. You need to correct the initial plan according to the feedback and output the refined plan.

User:

Goal: collect 3 stone.

Recipe: stone can be mined below the ground with at least wooden_pickaxe.

Initial Plan:

- 1. Mine 3 log with barehand
- 2. Craft 9 planks
- 3. Craft 2 stick
- 4. Craft 1 wooden_pickaxe
- 5. Dig down with wooden_pickaxe
- 6. Mine 3 log with wooden_pickaxe
- 7. Mine 3 stone with wooden_pickaxe

Feedback:

Incorrect. Step 4 "Crafting 1 wooden_pickaxe" must be done on a crafting_table. But the plan does not include "Craft 1 crafting_table" before Step 4, so the first mistake is "ignoring obtaining one crafting_table before using it to craft other tools". Step 6 "Mine 3 log with wooden_pickaxe" is executed after Step 5 "Dig down with wooden_pickaxe". But log only exists on above the ground. So the second mistake is "collecting materials on above the ground when the agent is already in below the ground".

Mistake:

1. Ignoring obtaining one crafting_table before using it to craft other tools.

2. Collecting materials on above the ground when the agent is already in below the ground.

Refined Plan:

- 1. Mine 3 log with barehand
- 2. Craft 9 planks
- 3. Craft 2 stick
- 4. Craft 1 crafting_table
- 5. Craft 1 wooden_pickaxe
- 6. Dig down with wooden_pickaxe
- 7. Mine 3 stone with wooden_pickaxe
- _____

Goal: {goal of new task} Recipe: {recipe of new task}

Initial Plan: {initial plan of new task}

Feedback: {feedback of the initial plan} Refined Plan:

Assistant:

Tree-of-thought (Yao et al., 2024). Tree-of-thought generates the plan for each task by exploring multiple possible paths simultaneously. Starting from the initial state, Tree-of-thought samples three possible steps each time that can be executed next. In our experiments, we adopt "OpenAI-gpt-4" as the backbone model. We set the temperature of LLM to 0.7 and sample for three times to get three possible next step each time. During planning, Tree-of-thought maintain multiple generated plans and select the best one through LLM rating. The generation of each plan ends when inferring the step to create the required item of the task goal.

The prompt used to sample possible steps is shown as follows:

System:

You are a helpful assistant in Minecraft. I will give you a goal to achieve in Minecraft, the recipe for the required item of the goal and the already generated partial plan. You need to generate the next step for current plan.

User:

Goal: collect 3 stone. Recipe: stone can be mined below the ground with at least wooden_pickaxe. Already Generated Plan: 1. Mine 3 log with barehand 2. Craft 9 planks 3. Craft 2 stick Next Step: 4. Craft 1 crafting_table ========

Goal: {goal of new task} Recipe: {recipe of new task} Already Generated Plan: {generated partial plan} Next Step:

Assistant:

After generating multiple candidate plans, the best plan is selected by taking the average result of three times of rating by LLMs. The prompt used to rate generated plan is shown as follows:

System:

Evaluate if the given plan can help to reach the goal (sure/likely/impossible)

User:

Plan: 1. Mine 3 log with barehand 2. Craft 9 planks 3. Craft 2 stick 4. craft 1 wooden_pickaxe Evaluation: impossible ========

Goal: {goal of new task} Plan: {generated plan of new task} Evaluation:

Assistant:

DEPS (Wang et al., 2023c). DEPS first generates the initial plan and then re-generates the plan through integrating the description and the explanation of the plan given by large language models. We adopt "OpenAI-gpt-4" as the backbone model.

DEPS first generates the initial plan by incontext learning in a code style:

System:

You are a helper agent in Minecraft. You need to generate the sequences of steps for a certain task in Minecraft.

User:

Goal: collect 3 stone.

Recipe: stone can be mined below the ground with at least wooden_pickaxe. The code for collecting 3 stone is as follows: def collect_3_stone(): mine({'log':3}, null); craft({'planks':9}, null); craft({'stick':2}, null); craft({'crafting_table':1}, null); craft({'wooden_pickaxe':1}, null); dig_down(wooden_pickaxe); mine({'stone':3}, wooden_pickaxe); return '3 stone'; ========

Goal: {goal of new task} Recipe: {recipe of new task}

Assistant:

Then, DEPS fixes the possible mistakes in the initial plan based on the description and the explanation of the initial plan:

System:

You are a helper agent in Minecraft. You need to generate the sequences of steps for a certain task in Minecraft.

User:

Goal: collect 3 stone. Recipe: stone can be mined below the ground with at least wooden_pickaxe. The code for collecting 3 stone is as follows: def collect_3_stone(): mine({'log':3}, null); craft({'planks':9}, null); craft({'stick':2}, null); craft({'wooden_pickaxe':1}, null); dig_down(wooden_pickaxe); mine({'stone':3}, wooden_pickaxe);

return '3 stone'; Descriptor: The step "craft({'wooden_pickaxe':1}, null)" will not be executed successfully. Explainer: Because I need a crafting table to craft the wooden pickaxe. RePlanner: The code for collecting 3 stone is as follows: def collect_3_stone(): mine({'log':3}, null); craft({'planks':9}, null); craft({'stick':2}, null); craft({'crafting_table':1}, null); craft({'wooden_pickaxe':1}, null); dig_down(wooden_pickaxe); mine({'stone':3}, wooden_pickaxe); return '3 stone'; _____

Goal: {goal of new task} Recipe: {recipe of new task} generated code for new task Assistant:

Plan-and-Solve (Wang et al., 2023b). Plan-and-Solve enables the large language model to generate the intermediate planning results before outputting the final plan. Plan-and-Solve can calculate the execution results of generated plans by in-context learning and thus enhances the accuracy of generated plans. We adopt "OpenAI-gpt-4" as the backbone model. The prompt used to generate the plans is shown as follows:

System:

You are a helpful assistant in Minecraft. I will give you a goal to achieve in Minecraft, the recipe for the required item of the goal. You need to generate the plan to achieve the task goal.

User:

Goal: collect 3 stone.

Recipe: stone can be mined below the ground with at least wooden_pickaxe.

Let's first understand the task and devise a plan to achieve the task goal.

Then, let's carry out the plan and achieve the task goal step by step.

Plan:

- 1. Mine 3 log with barehand
- 2. Craft 9 planks
- 3. Craft 2 stick
- 4. Craft 1 crafting_table
- 5. Craft 1 wooden_pickaxe
- 6. Dig down with wooden_pickaxe
- 7. Mine 3 stone with wooden_pickaxe

Task goal achieving:

1. I obtain 3 log by mining with barehand

2. I obtain 9 planks using 3 log

3. I obtain 2 stick using 2 planks, I have 7 planks left

4. I obtain 1 crafting_table using 4 planks, I have 3 planks left

5. I obtain 1 wooden_pickaxe using 2 stick and 3 planks on 1 crafting_table using, I have 0 planks left and 0 stick left

6. I reach below the ground by digging down with wooden_pickaxe

7. I obtain 3 stone by mining with wooden_pickaxe

In conclusion, I have achieved the task goal by obtaining 3 stone.

Goal: {goal of new task} Recipe: {recipe of new task}

Let's first understand the task and devise a plan to achieve the task goal.

Then, let's carry out the plan and achieve the task goal step by step.

Plan:

Assistant:

Openai-o1 (OpenAI, 2024). Openai-o1 uses the same prompt as Chain-of-thought. However, Openai-o1 will conduct more exploration on different paths and more refinements before outputting the final plans.

D.2 Baselines in Dynamic Planning

In dynamic planning experiments, the agent needs to interact with the environment and obtain the feedback to help conduct the planning. Therefore, the baseline models need to inform LLMs the feedback from the environment through prompts and utilize LLMs' reasoning ability to help generate proper plans.

React (Yao et al., 2023). React first conduct reasoning based on the task goal and the observation from the environment to infer the next few steps to be executed. Then it executes the inferred steps in the environment to obtain new observation, and repeats the reasoning and execution process until achieving the task goal. We adopt "OpenAI-gpt-4" as the backbone model. The used prompt is shown as follows:

System:

You are a helpful assistant in Minecraft. I will give you a goal to achieve in Minecraft, and you need to reason what steps to be executed next based on the task goal and the observation from the environment. Output the reasoning thought and the next few steps. After executing these steps in the environment, you will receive the observation from the environment. Then you need to reflect the execution results and continue to output the next few steps to be executed.

User:

Goal: collect 3 stone.

Recipe: stone can be mined below the ground with at least wooden_pickaxe.

Assistant:

Reasoning:

There is nothing in my inventory. I need to collect stone, so I first need to craft a wooden_pickaxe, arrive at below the ground, and mine the stone. Therefore, I need to craft the wooden_pickaxe. Steps: Mine 2 log with barehand Craft 8 planks Craft 2 stick Craft 1 crafting_table

Craft 1 wooden_pickaxe

User:

Observation:

Mine 2 log with barehand completed, there are 2 log in my inventory. Craft 8 planks completed, there are 0 log, 8 planks in my inventory. Craft 2 stick completed, there are 6 planks, 2 stick in my inventory. Craft 1 crafting_table completed, there are 2 planks, 2 stick, 1 crafting_table in my inventory. Craft 1 wooden_pickaxe failed.

Assistant:

Reasoning:

I failed crafting 1 wooden_pickaxe since there is no enough planks in my inventory. So I need to mine 1 more log to craft more planks. Then execute the following steps: Steps: Mine 1 log with barehand

Craft 1 planks

Craft 1 wooden_pickaxe

User:

Observation:

Mine 1 log with barehand completed, there are 1 log, 2 planks, 1 crafting_table, 2 stick in my inventory. Craft 1 planks completed, there are 0 log, 3 planks, 1 crafting_table, 2 stick in my inventory. Craft 1 wooden_pickaxe completed, there are 1 crafting_table, 1 wooden_pickaxe in my inventory.

Assistant:

Reasoning:

I have 1 wooden_pickaxe in my inventory now, I can dig down to arrive at below the ground to mine stone. Therefore I need to dig down with wooden_pickaxe and mine stone]

Steps:

Dig down with wooden_pickaxe Mine 3 stone with wooden_pickaxe

User:

Observation: Dig down with wooden_pickaxe completed, there are 1 crafting_table, 1 wooden_pickaxe in my inventory. I am already below the ground. Mine 3 stone with wooden_pickaxe completed, there are 1 crafting_table, 1 wooden_pickaxe, 3 stone in my inventory. Goal achieved.

User: Goal: {goal of new task} Recipe: {recipe of new task}

Assistant:

Jarvis-1 (Wang et al., 2023d). Jarvis-1 can adjust the plans based on the execution results from the environment. Through deep thinking, it can find the reason for the failure of previous plan and generate the corrected plan based on the thinking. We adopt "OpenAI-gpt-4" as the backbone model. The prompt used by Jarvis-1 is shown as follows:

System:

You are a helpful assistant in Minecraft. I will give you a goal to achieve in Minecraft and the initial plan to achieve the goal. After executing the plan in the environment, you will receive the observation about the execution results of the plan from the environment. You need to reflect the reason for the failure if the plan cannot be executed successfully and output corrected plan based on your reflection. Repeating the above process until the goal being achieved.

User:

Goal: collect 3 stone. Recipe: stone can be mined below the ground with at least wooden_pickaxe.

Initial Plan:

- 1. Mine 2 log with barehand
- 2. Craft 8 planks
- 3. Craft 2 stick
- 4. Craft 1 crafting_table
- 5. Craft 1 wooden_pickaxe
- 6. Mine 3 stone with wooden_pickaxe

Observation:

Mine 2 log with barehand completed, there are 2 log in my inventory. Craft 8 planks completed, there are 0 log, 8 planks in my inventory. Craft 2 stick completed, there are 6 planks, 2 stick in my inventory. Craft 1 crafting_table completed, there are 2 planks, 2 stick, 1 crafting_table in my inventory. Craft 1 wooden_pickaxe failed.

Assistant:

Thinking:

I failed crafting 1 wooden_pickaxe since there is no enough planks in my inventory. So I need to mine 1 more log to craft more planks.

Corrected Plan:

- 1. Mine 3 log with barehand
- 2. Craft 9 planks
- 3. Craft 2 stick
- 4. Craft 1 crafting_table
- 5. Craft 1 wooden_pickaxe
- 6. Mine 3 stone with wooden_pickaxe

User:

Observation:

Mine 3 log with barehand completed, there are 3 log in my inventory. Craft 9 planks completed, there are 0 log, 9 planks in my inventory. Craft 2 stick completed, there are 7 planks, 2 stick in my inventory. Craft 1 crafting_table completed, there are 3 planks, 2 stick, 1 crafting_table in my inventory. Craft 1 wooden_pickaxe completed, there are 1 crafting_table, 1 wooden_pickaxe in my inventory. Goal achieved.

User: Goal: {goal of new task} Recipe: {recipe of new task}

Initial Plan: {initial plan of new task}

Observation: {observation from the environment}

Assistant:



Figure 9: The performance of different models on plans of different lengths. Best viewed on screen.



Figure 10: Sample diversity of different models on tasks in different groups. Best viewed on screen.

E Supplementary Experiments

E.1 Advantage Analysis

To demonstrate the advantages of our proposed agent over other baselines, we evaluate the sensitivity of different models to plan length and the stability of inferring intermediate steps during planning. We first evaluate the planning performance of different models on plans of different lengths. As shown in Figure 9, the performance of baseline models varies significantly with the increase of the plan length, exhibiting extreme instability, especially for plans from the length of 7 to the length of 12. In contrast, our agent shows no significant variation in performance when the length of the plans increases. This is because our agent can narrow down the situations that the agent needs to consider for goal decomposition through the recursive goal decomposition module. Therefore, the difficulty of goal decomposition will not increase when the plans become longer, enabling our agent to conduct robust planning for various tasks.

Then we evaluate the stability of immediate step inference for different models. For tasks in each group, we provide a partial plan in the ground truth and require the model to generate the next three steps. For each task, we sample 50 generations for each model and calculate how many different three steps there are in these 50 generations, denoted as diversity value. As shown in Figure 10, compared to other baselines, our agent consistently maintains low diversity for the planning of different tasks. This demonstrates the stability of our agent in goal decomposition and plan generation.



Figure 11: Results of ablation study for state consistency maintaining, GPT4 as the backbone model, under the Accuracy metric.



Figure 12: Results of ablation study for state consistency maintaining, Llama-3-8B as the backbone model, under the Accuracy metric.



Figure 13: Results of ablation study for state consistency maintaining, GPT4 as the backbone model, under the Edit Distance metric.



Figure 14: Results of ablation study for state consistency maintaining, Llama-3-8B as the backbone model, under the Edit Distance metric.

E.2 Supplementary Ablation Study

We show part of the experimental results here that are not presented in the main part due to space limitations. For the ablation study in section 4.4, in addition to F1-Score, we also present the results under Accuracy and Edit Distance metrics. As shown in Figure 11 and Figure 12, under the Accuracy metric, our proposed state consistency maintaining module can significantly improve the quality of generated plans. Moreover, As shown in Figure 13 and Figure 14, under the Edit Distance metric, we can still draw the same conclusion because the proposed module can reduce the difference between generated plans and the ground truth. In addition, for the two methods to select anchor steps, the step scoring method can bring greater performance improvement, as the sliding window method chooses too many anchor step pairs and brings unnecessary noise to the plans.