

Cognitive Map for Language Models: Optimal Planning via Verbally Representing the World Model

Doyoung Kim¹ Jongwon Lee² Jinho Park¹ Minjoon Seo¹

¹KAIST AI ²Samsung Research
{doyoungkim, binlepain178, minjoon}@kaist.ac.kr, jay722.lee@samsung.com

Abstract

Language models have demonstrated impressive capabilities across various natural language processing tasks, yet they struggle with planning tasks requiring multi-step simulations. Inspired by human cognitive processes, this paper investigates the optimal planning power of language models that can construct a cognitive map of a given environment. Our experiments demonstrate that cognitive map significantly enhances the performance of both optimal and reachable planning generation ability in the Gridworld path planning task. We observe that our method showcases two key characteristics similar to human cognition: **generalization of its planning ability to extrapolated environments and rapid adaptation with limited training data.** We hope our findings in the Gridworld task provide insights into modeling human cognitive processes in language models, potentially leading to the development of more advanced and robust systems that better resemble human cognition.

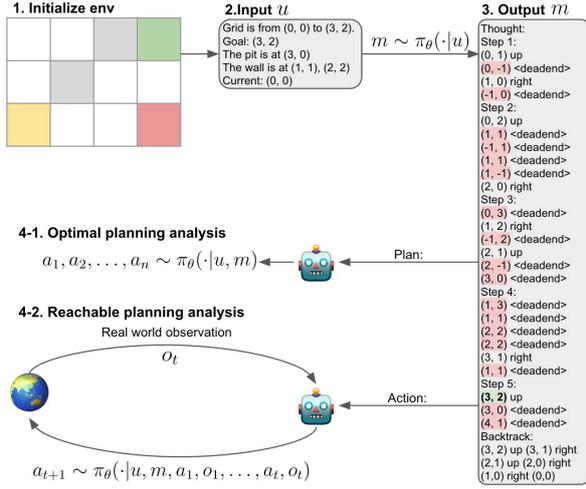
1 Introduction

Language models have recently demonstrated remarkable proficiency in a variety of complex tasks (Brown et al., 2020; Touvron et al., 2023; Chowdhery et al., 2023; Chen et al., 2021), from natural language understanding to code generation, primarily through the training objective of next-token prediction. This training paradigm enables language models to excel at general planning tasks by leveraging their extensive learned knowledge and pattern recognition capabilities (Ahn et al., 2022; Liang et al., 2023; Song et al., 2023). However, it is also true that they often falter in scenarios that require robust, long-horizon planning tasks (Dziri et al., 2024). It is in contrast that humans naturally employ model-based planning, internally construct models to simulate outcomes and guide optimal decision-making, as extensively documented in cognitive science literature (Daw et al.,

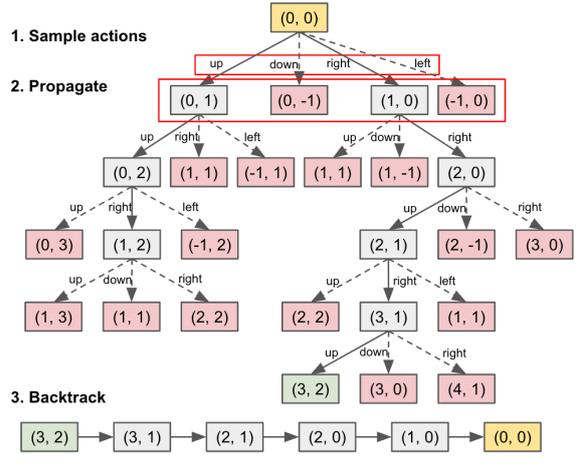
2005). This distinction aligns with the dual-process theory of reasoning: System 1 processes are fast, automatic, and pattern-based, akin to model-free planning, while System 2 processes are slower, deliberative, and involve explicit reasoning (Daniel, 2017).

Existing methods suggest different ways to model System 2 process. Exploration-based planning methods such as Tree of Thought (ToT) (Yao et al., 2023) use external simulation via multi-turn inference to enhance the exploration ability of language models. While demonstrating strong exploration capabilities, they fall short in achieving optimal planning because they conduct explicit rollouts rather than incorporating the simulation process within the action decisions. On the other hand, the Chain of Thought (CoT) (Wei et al., 2023) methodology is also widely used to model the System 2 process of language models, injecting an optimal path into the reasoning process. Despite its impressive enhancement in reasoning and planning, existing CoT methods focus on generating a solution from start to goal, which does not align with the true notion of "simulation," i.e., analyzing real-world systems and predicting outcomes. Naive CoT only provides intermediate steps from the start to the goal, which is inappropriate as a simulation process.

When humans solve planning tasks like Gridworld, they typically conduct iterative simulations until reaching the goal, while avoiding deadend states. After reaching the goal, they work backward to figure out the current move to get to the goal state. This simulation does not require actual rollout but rather the construction of a "cognitive map." Based on the intuition, we train language models with datasets augmented with a cognitive map, comprising three key processes: **Sampling**, **Propagation**, and **Backtracking**. Namely, the model samples plausible actions that can be performed at each state. It then propagates for each



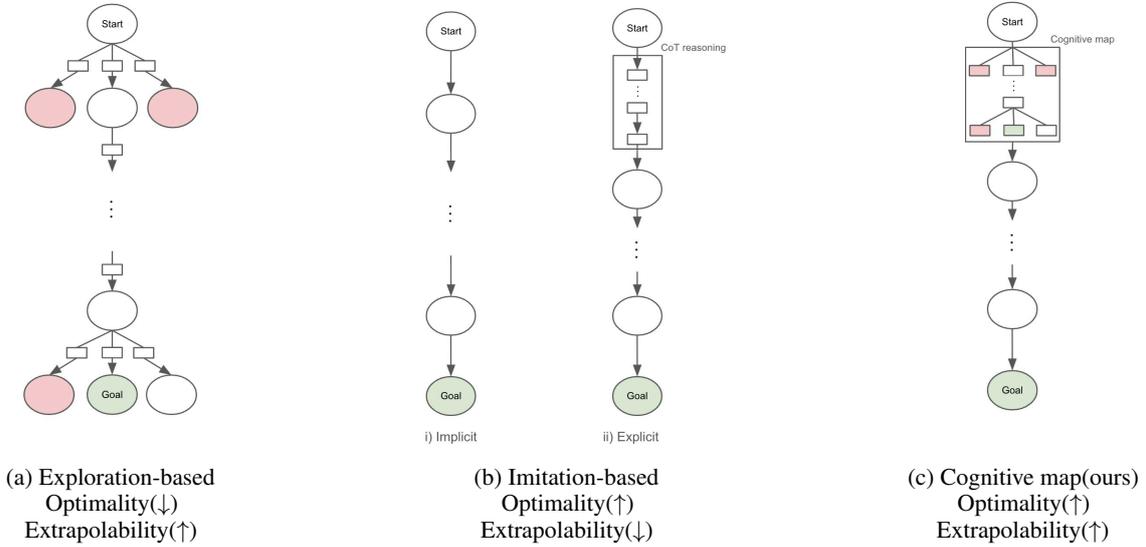
(a) Overall process



(b) Structure of the cognitive map. The construction is done by the generation of π_θ without any actual interaction.

Figure 1: Overview of the paper(Figure 1a): 1. We initialize the world environment which the model is going to interact. 2. The textualized input instruction containing information about the environment(u) is fed into the model. 3. Before interacting with the world, the model constructs a textualized cognitive map(m). Our construction of m is a tree-structured verbal representation of the world model(Figure 1b), which is represented in a sequential manner(We provide details for sampling, propagation, and backtrack in Section 4). 4. With the constructed map, the agent interacts with the environment. We analyze the power of cognitive map in generating both the optimal plan without further observation(4-1) and the reachable plan with partial observation(4-2).

We show that the cognitive map deduces optimal plan, and it shows human-cognitive characteristics such as generalization to extrapolated environments or rapid adaptation with limited training data(Section 6).



(a) Exploration-based
Optimality(↓)
Extrapolability(↑)

(b) Imitation-based
Optimality(↑)
Extrapolability(↓)

(c) Cognitive map(ours)
Optimality(↑)
Extrapolability(↑)

Figure 2: Different approaches of planning sequence of actions(\rightarrow) from the start state(Start) to the goal state(Goal). Figure 2a: Exploration-based planning (Yao et al., 2023) samples plausible actions for each state, creates an explicit tree, and explores with a tree search algorithm while avoiding deadend state(red circle)(there is an optional reasoning(\square) before taking an action(\rightarrow)). Exploration-based planning focuses on making a reachable plan while sacrificing the optimality of the plan. Figure 2b: Imitation-based planning (Yao et al., 2022b) selects actions(\rightarrow) based on the episodic trajectory. It could either directly select actions for each state without further reasoning(b-i, implicit), or use Chain of Thought reasoning (Wei et al., 2023)(\square) of the decision(b-ii, explicit). Imitation-based planning can successfully generate optimal plans within the interpolated data while struggling with extrapolated data. Figure 2c: Cognitive map construction is an extension of imitation-based planning, where the model first constructs a cognitive map that interprets the relations of global states before taking an action(\rightarrow). When the construction is done, it performs actions(\rightarrow) based on the constructed map. We claim that constructing cognitive maps ensures the optimality of the plan even in extrapolated data points.

action that reaches to a new state. The model explores(samples and propagates) potential actions and outcome states until the goal is reached. Once the goal is reached, it backtracks traces back from the goal state to the start state to refine the optimal path(See Section 4 for details).

We show that models fine-tuned using this approach exhibit significant characteristics similar to human cognitive map, which are:

- **Extrapolation:** The ability to solve problems in larger environments which is unseen during the training stage.
- **Rapid Adaptation:** Learning effectively with a small training dataset.

In this paper, we aim to show the strength of constructing cognitive maps in path planning Gridworld¹. By leveraging the global world representations that the cognitive map gives to the language model, we aim to model the cognitive process of humans, thereby enhancing the performance of language models in complex, long-horizon planning tasks.

2 Related work

2.1 Planning in cognitive science

Dual-process theories of cognition distinguish between System 1 and System 2 reasoning, where System 1 involves fast, automatic, and intuitive thinking, and System 2 involves slow, deliberate, and analytical thought (Daniel, 2017). System 1 and System 2 reasoning aligns with model-free and model-based planning, respectively. Unlike model-free planning, model-based planning benefits by constructing and utilizing internal models to simulate outcomes and guide decisions (Daw et al., 2005). Evidence from cognitive science suggests that humans rely on model-based planning for complex decision-making tasks, using internal simulations to evaluate the consequences of different actions (Doll et al., 2012).

One of the key characteristics of human cognition is extrapolation. Human ability to solve problems in larger, more complex environments that were not encountered during the initial learning phase has been largely observed across cognitive science domains (Yousefzadeh and Mollick, 2021;

¹We choose Gridworld since the world is easily scalable, so we can easily test the planning ability for extrapolated data. We further discuss the application of our method in general agent tasks in Section 7

Fellini and Morellini, 2011; Stojic et al., 2018). This capability allows humans to apply learned knowledge to novel situations, demonstrating a remarkable level of generalization. Additionally, human cognition is characterized by rapid adaptation, the ability to learn effectively from a small amount of data (Lake et al., 2017). This enables humans to quickly adjust their strategies and behaviors based on limited experiences or feedback, which is essential for navigating dynamic and unpredictable environments. These traits underscore the flexibility and efficiency of human cognitive processes, setting a benchmark for artificial intelligence systems aiming to replicate human-like problem-solving and learning abilities.

2.2 Planning in language model domain

Language models are primarily trained using next-token prediction, which facilitates model-free planning by allowing them to generate text based on learned patterns and associations from vast amounts of data. This training approach enables language models to perform zero-shot planning in various domains, such as general planning and code generation, without requiring task-specific training (Chen et al., 2021). Language models excel at pattern matching and leverage their learned world model to plan and generate coherent text in given tasks (Brown et al., 2020).

Despite their successes, language models often struggle with planning tasks that require extensive foresight and detailed reasoning. Studies have shown that language models perform poorly on tasks requiring robust planning and complex decision-making, highlighting the limitations of model-free approaches (Xie et al., 2024; Valmeekam et al., 2022).

2.3 Exploration-based planning

One limitation of language model planning is the lack of opportunity to explore a wide range of world states. To address this, exploration-based methods have been designed to enable language models to explore more effectively. One such example is Tree of Thought (ToT) (Yao et al., 2023), which leverages language models' knowledge to sample child nodes and simultaneously evaluate the value of the current state to structure and search the tree structure at the same time. Also one can enhance the searching by using Monte Carlo Tree Search (MCTS) algorithm (Kocsis and Szepesvári, 2006) to explore states and update their value (Hao

et al., 2023; Zhou et al., 2023). These exploration-based planning methods enhance the model’s ability to reach the goal, with a perfect probability if we have infinite time and inference cost. However, their goal is to enhance the reachability of the plan while sacrificing the optimality of the generated plan.

2.4 Imitation-based planning

One way to inject optimal planning capability is by imitating the optimal demonstrations. We call it as Chain of Thought (CoT)(Wei et al., 2023) methodology, and it has demonstrated significant improvements in language models’ reasoning and planning abilities. They involve guiding the model through a series of chains generated by human or ground truth that simulates step-by-step reasoning, enhancing the model’s decision-making processes via few-shot demonstration or fine-tuning. They imitate ground truth reasoning demonstrations to improve the performance in tasks requiring complex reasoning (Nye et al., 2021; Yao et al., 2022b).

One strength of training with CoT is that the demonstration does not have to be optimal. Methods such as Searchformer (Lehnert et al., 2024) experimentally showed that imitating automatic tree search planning such as A* rather than human demonstration is also helpful when planning. The searching algorithm conducts tree search via heuristic cost function, so it can cover more world model representation than ground truth demonstration. Yet they sacrifice the optimality of the plan, they show that the policy of the model can further imply more optimal solutions than the demonstrated algorithm.

Limitation of imitation-based planning in extrapolated data Despite the great potential of expressiveness power of CoT, imitation-based planning does not work well, even for frontier language models such as GPT 4 (OpenAI et al., 2024). They are shown to struggle at solving and generalizing compositional tasks such as multiplication, Einstein’s puzzle (Dziri et al., 2024) or Blocksworld puzzle (Valmeekam et al., 2022; Stechly et al., 2024) in an extrapolated data even after fine-tuning or giving few-shot demonstrations. It is in contrast with features of human cognition planning.

One of the biggest reasons for the fallacy is that human policy is compact and entangled so they do not let language models see the global representation of the world model. Such entanglement often

lets hallucination happen during generation (Peng et al., 2024). Theoretical approaches highlight the effectiveness of CoT (Merrill and Sabharwal, 2024; Feng et al., 2024), but they also note that CoT demonstration should be long enough to potentially solve planning problems that necessitate System 2 reasoning.

To sum up, exploration-based planning method lacks in generating the optimal plan, while imitation-based planning lacks in succeeding in the extrapolated data(Figure 2). We claim that imitating cognitive map deduces optimal model-based planning in extrapolated data points. By explicitly reasoning the representation of global states, we hypothesize that language models can achieve optimal planning capabilities. This approach mirrors human cognitive processes and enables language models to handle complex, long-horizon planning tasks more effectively (Daw et al., 2005).

3 Planning task interacting with world model

Planning task with environment feedback in language model can be formalized as a Markov decision process with instruction space \mathcal{U} , state space \mathcal{S} , action space \mathcal{A} , observation space \mathcal{O} , metric space \mathcal{C} , transition function $T : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$, and metric function $R : \mathcal{U} \times (\mathcal{S} \times \mathcal{A})^n \rightarrow \mathcal{C}$ where n is the trajectory length. Note that for language model domain, \mathcal{U} , \mathcal{A} , and \mathcal{O} are given as sequences of language tokens.

Given an instruction $u \in \mathcal{U}$, the model θ first generates the action $a_1 \sim \pi_\theta(\cdot|u) \in \mathcal{A}$ according to its policy π_θ . For each state $s_t \in \mathcal{S}$ and its observation $o_t \in \mathcal{O}$, the agent generates the corresponding action in the $t + 1$ step $a_{t+1} \sim \pi_\theta(\cdot|u, a_1, o_1, \dots, a_t, o_t) \in \mathcal{A}$, which concludes to a new state $s_{t+1} = T(s_t, a_t)$ and its observation o_{t+1} . The interaction loop repeats until the task has terminated for some reason(succeeded, failed, number of steps exceeded maximum value, etc.), and the action trajectory is denoted as:

$$e = (u, a_1, o_1, \dots, o_{n-1}, a_n, o_n) \sim \pi_\theta(e|u),$$

$$\pi_\theta(e|u) = \prod_{j=1}^n \pi_\theta(a_j|u, a_1, o_1, \dots, o_j),$$

Finally, we define the (action, state) trajectory $k = ((a_1, s_1), \dots, (a_n, s_n))$ accordingly. The final reward is computed based on the metric function $R(u, k)$.

Note that we can apply reasoning before deciding the action with so-called a "thinking" process. Namely, we have a thinking space \mathcal{M} (of language subset) which generates $m \sim \pi_\theta(\cdot|u) \in \mathcal{M}$, then generate $a_1 \sim \pi_\theta(\cdot|u, m) \in \mathcal{M}$ upon the generated thought.

4 Cognitive map of a given world model

Algorithm 1 Cognitive map construction algorithm

```

map  $\leftarrow$  [], queue  $\leftarrow$  [start]
while (goal,  $\cdot$ )  $\notin$  map do
  s  $\leftarrow$  POP(queue)
  for a  $\in$  S(s) do ▷ Sampling
    child  $\leftarrow$  T(s, a) ▷ Propagation
    if child  $\notin$  DEADEND then
      map+ = (child, a)
      queue+ = child
    end if
  end for
end while
s  $\leftarrow$  goal, backtrack  $\leftarrow$  []
while start  $\notin$  backtrack do
  backtrack+ = s
  (s, a)  $\leftarrow$  GET((s,  $\cdot$ ), map)
  s  $\leftarrow$  T-1(s, a) ▷ Backtrack
end while

```

Our proposed process can be divided into three different stages: sampling, propagation, and backtrack. The sampling stage can be defined as the process where the model expects potential actions that can be applied for each state. The propagation stage explores potential outcomes of the actions sampled from the sampling stage. The backtrack stage is the process of tracing back from the goal state to refine and select the optimal path based on the outcomes of the simulation (See Algorithm 1 for the pseudocode).

4.1 Sampling

During the sampling stage, the model samples possible actions in $S(s) \subset A$ regarding the current state s , which leads to a new state that was not propagated before. This sampling process is repeated until reaching the desirable goal. This iterative approach enables the cognitive map of the agent to explore the various states of the given world model.²

²In this paper, we treat only for Gridworld case so there are small number of action candidates to sample (there are always

4.2 Propagation

During the propagation stage, the model expects a transited state $T(s, a) \in \mathcal{S}$ for each action a sampled from the current state s . For example, in the Gridworld, the model simulates different paths by considering movements in four directions (up, down, left, right) from each cell. The goal is to explore various routes to reach the target cell, considering obstacles and the grid's boundaries. This stage is crucial for understanding the global representation of the world model by expecting the future consequences of each action before committing to the actual decision.

4.3 Backtrack

Once the sampling and propagation stages have concluded by reaching a desirable goal state, we need to backtrack through the simulated paths to determine the most efficient route taken to achieve the goal. This involves assessing the propagated paths and identifying the one that reaches the goal. Backtrack search identifies $T^{-1}(s, a) \in \mathcal{S}$ for each state-action pair to ensure that the selected path is valid. By refining the decisions made during the simulation stage, backtracking ultimately guides the model to make informed, strategic choices based on the simulated outcomes. Since the propagation stage deduces only one path with minimal steps, we can guarantee the optimality of the generated path.

4.4 Training how to generate cognitive map

To sum up, constructing the cognitive map m is a sequential application of sampling(S), propagation(T), and backtrack(T^{-1}). We train the language model θ via supervised learning method so that the model can successfully construct the cognitive map without any external interaction, as shown in Figure 1b.

5 Experimental setup

5.1 Basic setup

Task, train, and inference detail In this paper, we set textualized Gridworld (Brown, 2015) as the main task. Gridworld is a task that involves path planning from the start state to the goal state while avoiding the pit, wall, and grids outside the world. Especially, we ensure that there is only one path

four moves, up, down, left, or right). We discuss possible sampling strategies that the agent can make in a bigger world in Section 7.

from start to goal for each environment (For the example of input including instruction and the world information fed to the model, refer Appendix A.1).

We train the model for 1 epoch with 50K samples of size 10×10 at largest. After training, we test the model with 3K samples with each grid being 20×20 at largest. We utilize Llama-3-8B model³ throughout the whole experiments. In each turn, we set the maximum token length of the model to be 8192. (For more details about experimental setup, refer Appendix A.2).

Optimal planning analysis: single-turn setting

We see the capability of cognitive map on generating the optimal plan without any further observations (Corresponding to 4-1 of Figure 1a). We evaluate the optimality of the trajectory generated by $a_1, a_2, \dots, a_n \sim \pi_\theta(\cdot | u, m) \in \mathcal{A}^n$. This is equivalent to generating the whole plan in a single-turn, so we train the model to generate the whole plan in one turn and evaluate it.

Given the instruction u and generated (action, state) trajectory k , we first define optimal path k^* as an element of $\arg \min_{k \in \{k | r(u, k) = \text{success}\}} |k|$. Since we ensure there is only one possible k^* , we only need to check whether the generated trajectory is k^* . Now we can define $R(u, k)$ as follows:

- $R(u, k) = \text{SUCCESS}$ if $k = k^*$
- $R(u, k) = \text{FAIL}$ otherwise

Reachable planning analysis: multi-turn setting

We also investigate whether the cognitive map can deduce a reachable plan (which doesn't require optimal planning) in a multi-turn setting (Corresponding to 4-2 of Figure 1a). We evaluate the reachability of the trajectory a_1, a_2, \dots, a_n generated by $a_{t+1} \sim \pi_\theta(\cdot | u, m, a_1, o_1, \dots, a_t, o_t) \in \mathcal{A}$. There are multiple interactions between the model and the environment in a multi-turn setting, so failure cases can be divided into three cases (deadend, max step, and invalid). Especially, given the instruction u and generated (action, state) trajectory k , we define $R(u, k)$ as follows:

- $R(u, k) = \text{SUCCESS}$ if $k[-1][1] = \text{goal}$
- $R(u, k) = \text{DEADEND}$ if $k[-1][1] \in P$
- $R(u, k) = \text{MAX STEP}$ if $|k| > \text{max}$

³<https://llama.meta.com/llama3/>

- $R(u, k) = \text{INVALID}$ if $\exists a \in k[0] | a \notin A$

⁴ In the paper, we set the maximum steps $\text{max} = 200$ for each environment. We denote P as the set of deadend states, i.e. set of pits and walls in the grid, and all the states outside the grid.

For each observation, we give the information of the current state along with possible moves that do not directly result the deadend state along with its corresponding states. For example, if the current state s_t is (11, 4) and there the possible actions are right or left, the observation o_t is "Current:\n(11, 4)\nPossible:\n(10, 4)\nleft\n(12, 4)\nright".

5.2 Cognitive map design choices

Our construction of the cognitive map is straightforward; We sample all 4 possible directions ("up", "down", "right", and "left") as actions, and we propagate each action iteratively until reaching the goal state. After reaching the goal, we backtrack until reaching the start state.

Marking deadend In this paper, we variate two different cognitive maps when marking the deadend state. We can either verbalize all samples and mark the actions resulting deadend state with a special token (`<deadend>`) (denoted as "MARKING deadend"), or just verbalize all the actions including the deadend states (denoted as "w.o. MARKING"). For example, the verbalized cognitive map at state (0, 0) in Figure 1b is as follows:

MARKING deadend: "(0, 1) up (0, -1) `<deadend>` (1, 0) right (-1, 0) `<deadend>`"
w.o. MARKING: "(0, 1) up (0, -1) down (1, 0) right (-1, 0) left".

See Appendix B.1 for the whole construction.

Baselines We name NONE and COT as our baseline experiments representing imitation-based planning. NONE implicitly learns how to conduct path planning, while COT learns a verbalized backward trace as a Chain of Thought (CoT) demonstration. For example, the COT reasoning of the Figure 1b is "(3, 2)up\n(3, 1)\nright\n(2, 1)\nup\n(2, 0)\nright\n(1, 0)\nright\n(0, 0)", while we do not verbalize anything for NONE. See Appendix B.1 for the whole construction.

Backward cognitive map construction: BWD

Backward chaining is a powerful approach that simplifies complex problems by focusing on the

⁴Although our main analysis consists only of success rate, we provide plots for different types of fails in Appendix C.3

desired outcome and systematically working back to the starting point. LAMBADA (Kazemi et al., 2023) shows that backward chaining helps in reasoning tasks. We adopt the intuition to see if constructing the cognitive map in a backward manner enhances the planning ability of the language model.

In this paper, we define two types of construction, FWD and BWD. For FWD, the construction from the start to the goal is identical to the procedure stated in Section 4. For BWD, we build the reversed cognitive map starting from the goal state to the start state. We have the sampling function identical to $FWD(S)$. For the propagation stage, we expect a reverse transition state $T^{-1}(s, a) \in \mathcal{S}$ for a given state s and action a . For the backtracking stage, we backtrack the path from start to goal by iteratively searching $T(s, a) \in \mathcal{S}$. See Appendix B.2 for the example.

Additional setups We also experiment effects of excluding each component when designing the cognitive map. First, we see the effect of excluding the backtrack stage (denoted as "w.o. BACKTRACK"). Also, we observe the effect when we only sample possible moves instead of all moves (denoted as "w.o. ALL"). See Appendix B.1, B.2 for examples.

5.3 Comparison of the cognitive model with exploration-based planning

We also compare the planning ability of our method with exploration-based planning. We train the language model which traverses the world through a tree search algorithm via Depth-First Search (DFS), denoted as DFS (Pseudocode of the algorithm is in Algorithm 2). See Table 9 for the example.

6 Results and analysis

6.1 Overall result analysis

Cognitive map compared with baselines Table 1 shows both the optimal and reachable rate of plans generated by each experiment. As shown in the table, w.o. MARKING shows the best performance among the experiments both for BWD and FWD cognitive map construction in the optimal planning. (76.5% and 61.8% for BWD and FWD construction, respectively). Unlike in the optimal planning setting, MARKING deadend shows the best performance both for BWD and FWD cognitive map construction (88.5% and 85.4% for BWD and FWD construction, respectively). **Our experiments**

show that cognitive maps improve performance in the Gridworld path planning task. It boosts the optimal planning performance by up to 57.5% and reachable planning by up to 56.4% with implicit fine-tuning baseline, and enhancing optimal planning by up to 50% and reachable planning by up to 54.6% with CoT fine-tuning baseline. We also qualitatively show the planning performance in both optimal and reachable settings with respect to the size of the grid in Table 2.

Cognitive map further enhances reachability

A notable observation is the performance gap between optimal and reachable planning. Reachable planning experiments show significant improvements compared to optimal planning across all configurations, except for ALL BACKTRACK in FWD cognitive map construction. For instance, w.o. BACKTRACK from FWD map construction more than doubled its score (42.3% to 85.2%). This implies that **although the cognitive map was designed to find the optimal plan, it also substantially enhances reachable planning capability.**

BWD approach enhances the performance of the cognitive map

As shown in Table 1, both the optimal and reachable planning got the highest performance with BWD cognitive map construction. This aligns with the findings in LAMBADA (Kazemi et al., 2023), where they show that backward chaining helps in reasoning tasks.

6.2 Extrapolation ability

As stated in Section 5.1, we train the model on world sizes of up to 10x10 and test it on world sizes of up to 20x20 world to investigate the extrapolation ability. Extrapolation success is defined as the ability to succeed in rollouts beyond the 10x10 boundary. As shown in Table 2, the darkness outside the boundary of the red box shows that **the cognitive map helps in planning on the extrapolated data.** We also observe a consistent tendency for experiments with ALL inclusion, which have a wider coverage of success rate in the extrapolated data. These results align with the result discussed in Section 6.1. We qualitatively analyze all the results in Appendix C.2 and C.3.

6.3 Rapid adaptation

We experiment speed of the language model learning robust BWD MARKING deadend cognitive map construction by both watching its reachability and loss curve. As shown in Figure 3a, the success rate

	Implicit baseline	Explicit baseline	Cognitive map	
Optimal	NONE	CoT	MARKING deadend	w.o. MARKING
BWD	0.190	0.265	0.705	0.765
FWD	0.190	0.252	0.585	0.618
Reachable	NONE	CoT	MARKING deadend	w.o. MARKING
BWD	0.321	0.287	0.885	0.724
FWD	0.321	0.339	0.854	0.816

Table 1: Optimal and reachable rate of generated plans via single- and multi-turn settings: The first two columns(NONE and BACKTRACK) are the baselines for imitation-based learning, and the rest are different design choices of constructing the cognitive map. Also BWD constructs the map starting from the goal state, while FWD starts from the start state. See Appendix B for actual prompts.

of the model converges at the early stage(79.13% at step 500, 74.5% at step 625), which implies that it is easy for a model to learn how to generate a robust cognitive map. We also observe that the loss curve of the model converges to 0 much faster than other baseline methods(Figure 3b). **It suggests that the language model rapidly learns how to construct the cognitive map.**

6.4 Compare with exploration-based planning

Despite the significant planning ability demonstrated by the cognitive map, we observe that there is still room for improvement, particularly in reachable planning(See Appendix D for detailed setting and qualitative analysis). The results show that exploration-based planning is more effective in generating reachable plans(See Table 3). We interpret the incident due to an inductive bias of exploration-based planning toward exploring via tree search algorithm. With an infinite amount of time, tree search algorithm almost perfectly ensures the reachability toward the goal.

On the other hand, we show that constructing the cognitive map has benefits in generating the plan with optimal steps(See Figure 4). We see that the exploration-based planning has a mean inference steps for reaching the goal for n optimal steps as $O(n^2)$, while our method shows the optimal performance also for the reachable plan generation.

7 Discussion

Why is cognitive map easier to learn than CoT or implicit policy? One of the left curiosity is interpreting the result we obtained from Section 6.3, which implies that the language model rapidly learns how to construct the cognitive map. The re-

sult suggests that the language model easily learns how to construct a cognitive map. How is it possible?

In this paper, the generation of cognitive map $m \in \mathcal{M}$ given an instruction $u \in \mathcal{U}$ is modeled as a composition of following functions: sampling(S), propagation(T), and backtrack(T^{-1}). Given the current state s_t , the sampling strategy selects actions $S(s_t) \subset \mathcal{A}$. For each action $a \in S(s_t)$, propagation generates the next state $T(s_t, a)$. These steps iterate until we reach the goal state. Once the goal state g is reached, we backtrack to the start state $T^{-1}(s_t, a)$. Assuming that applying the transition function T and its inverse T^{-1} is trivial for the model, each generation step of m is straightforward. Hence if conducting each function is straightforward, we can also say that the overall generation of $m \sim \pi_\theta(\cdot|u) \in \mathcal{M}$ is also trivial.

Why does cognitive map deduce robust planning? Human demonstration of path planning is often abstract and entangled. This entanglement occurs between the human policy π_{human} and the replication policy from observed human demonstrations $\pi_{\text{human demo}}$. $\pi_{\text{human demo}}$ tends to be more complex and interwoven since it has various hidden factors. Namely, our objective is to expect

$$e = (u, a_1, o_1, \dots, o_{n-1}, a_N) \sim \pi_{\text{human}}(e|u),$$

but the human demonstration is given as

$$e' = (u, a'_1, o'_1, \dots, o'_{n-1}, a'_n) \sim \pi_{\text{human demo}}(e'|u),$$

where $(a'_1, \dots, a'_n) \subset (a_1, \dots, a_N)$, $n \ll N$.

Our underlying assumption is that π_{human} is relatively simple, and we can model it by composing sampling(S), propagation(T), and backtrack(T^{-1})

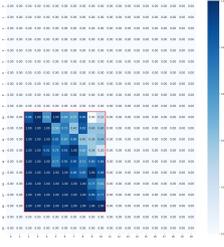
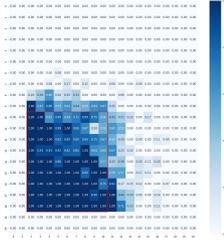
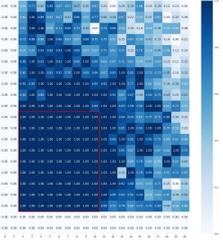
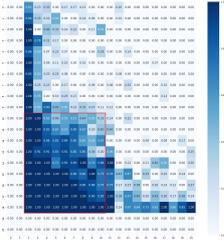
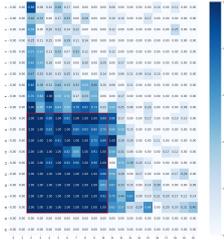
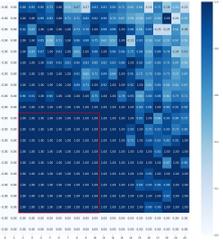
	Implicit baseline	Explicit baseline(best)	Cognitive map(best)
Optimal	 NONE: 19.0%	 BWD BACKTRACK: 26.5%	 BWD w.o. MARKING: 76.5%
Reachable	 NONE: 32.1%	 FWD BACKTRACK: 33.9%	 BWD MARKING deadend: 88.5%

Table 2: Qualitative comparison of reachable and optimal plan generation rate between baselines and our methods. The degree of the darkness at (x, y) coordinate of each plot denotes the performance of the corresponding model in Gridworld of size $x \times y$. The red box denotes the boundary of the training data. We provide visualization of all the experiments in Appendix C.

Cognitive map(best)	Exploration-based
0.885	0.945

Table 3: Comparison of reachability with exploration-based planning(DFS) and our best(BWD MARKING deadend cognitive map)

functions. We saw in the Gridworld that these functions are straightforward and easily replicated. By employing these simple and well-defined functions, the model effectively addresses the problem of data path planning. This approach allows the model to generate robust planning ability, including extrapolation beyond the training data.

Extension to general world Our discussion of planning is restricted to the Gridworld environment throughout this paper. Gridworld has a relatively small world size compared to general agent tasks such as those performed by web agents (Yao et al., 2022a) or travel agents (Xie et al., 2024). To extend planning capabilities to these larger, more complex worlds, we may need to modify the sampling strategy. The current sampling approach may not scale well in larger environments. A potential solution is

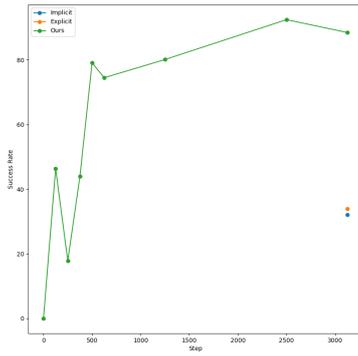
to train the model to sample the top n preferable actions, which would allow the model to efficiently navigate and plan in much larger and more intricate worlds. We leave the investigation of sampling strategy as a further research direction.

8 Acknowledgement

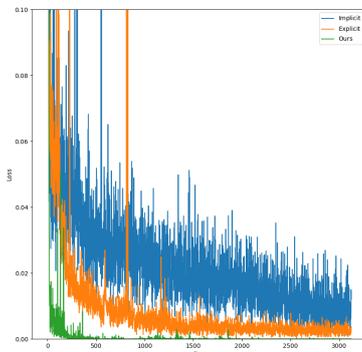
We thank Seonghyeon Ye, Hyeonbin Hwang, Se-june Joo and Juyoung Suk for their invaluable advice and support throughout this research.

References

Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Daniel Ho, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Eric Jang, Rosario Jauregui Ruano, Kyle Jeffrey, Sally Jesmonth, Nikhil J Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Kuang-Huei Lee, Sergey Levine, Yao Lu, Linda Luu, Carolina Parada, Peter Pastor, Jornell Quiambao, Kanishka Rao, Jarek Rettinghouse, Diego Reyes, Pierre Sermanet, Nicolas Sievers, Clayton Tan, Alexander Toshev, Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Mengyuan Yan, and Andy Zeng. 2022. *Do as i can,*



(a) Success rate comparison: For our method, we provide additional performance of the model checkpoints at step 125, 250, 375, 500, 625, 1250, 2500, and 3125.



(b) Loss curve comparison

Figure 3: Success rate and loss curve comparison among Implicit baseline(NONE: blue), Explicit CoT baseline(FWD BACKTRACK: orange), and our method(BWD MARKING deadend: green).

not as i say: Grounding language in robotic affordances.

Brandon Brown. 2015. *Q-learning with neural networks: Learning gridworld with q-learning*. Accessed on June 19, 2024.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul

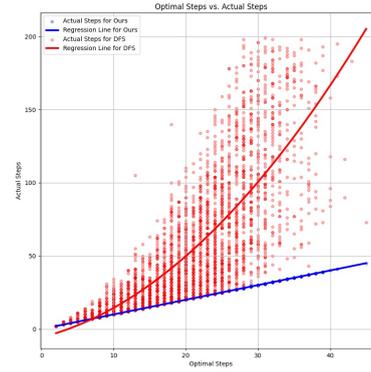


Figure 4: Number of inference steps per optimal steps for succeeded instances. We compare exploration-based planning(DFS: red) with our method(BWD MARKING deadend: blue) in the reachable planning setting. Each scattered dot denotes the succeeded instance, and the line plot denotes the regression line of the success. Note that the regression line of our method is almost identical to a $x = y$ graph, which means that cognitive map helps generating optimal performance even for the reachable planning setting.

Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.

Kahneman Daniel. 2017. *Thinking, fast and slow*.

Nathaniel D Daw, Yael Niv, and Peter Dayan. 2005. Uncertainty-based competition between prefrontal and dorsolateral striatal systems for behavioral control. *Nature neuroscience*, 8(12):1704–1711.

Bradley B Doll, Dylan A Simon, and Nathaniel D Daw. 2012. The ubiquity of model-based reinforcement learning. *Current opinion in neurobiology*, 22(6):1075–1081.

Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jiang, Bill Yuchen Lin, Sean Welleck, Peter West, Chandra Bhagavatula, Ronan Le Bras, et al. 2024. Faith and fate: Limits of transformers on compositionality. *Advances in Neural Information Processing Systems*, 36.

Laetitia Fellini and Fabio Morellini. 2011. Geometric information is required for allothetic navigation in mice. *Behavioural brain research*, 222(2):380–384.

Guhao Feng, Bohang Zhang, Yuntian Gu, Haotian Ye, Di He, and Liwei Wang. 2024. Towards revealing the mystery behind chain of thought: a theoretical perspective. *Advances in Neural Information Processing Systems*, 36.

Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. 2023.

- Reasoning with language model is planning with world model. *arXiv preprint arXiv:2305.14992*.
- Mehran Kazemi, Najoung Kim, Deepti Bhatia, Xin Xu, and Deepak Ramachandran. 2023. [Lambda: Backward chaining for automated reasoning in natural language](#).
- Levente Kocsis and Csaba Szepesvári. 2006. Bandit based monte-carlo planning. In *European conference on machine learning*, pages 282–293. Springer.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. [Efficient memory management for large language model serving with pagedattention](#).
- Brenden M. Lake, Tomer D. Ullman, Joshua B. Tenenbaum, and Samuel J. Gershman. 2017. [Building machines that learn and think like people](#). *Behavioral and Brain Sciences*, 40:e253.
- Lucas Lehnert, Sainbayar Sukhbaatar, Paul Mcvay, Michael Rabbat, and Yuandong Tian. 2024. Beyond a*: Better planning with transformers via search dynamics bootstrapping. *arXiv preprint arXiv:2402.14083*.
- Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. 2023. [Code as policies: Language model programs for embodied control](#).
- Ilya Loshchilov and Frank Hutter. 2017. [Sgdr: Stochastic gradient descent with warm restarts](#).
- William Merrill and Ashish Sabharwal. 2024. [The expressive power of transformers with chain of thought](#).
- Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, et al. 2021. Show your work: Scratchpads for intermediate computation with language models. *arXiv preprint arXiv:2112.00114*.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kopic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael

- Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. [Gpt-4 technical report](#).
- Binghui Peng, Sridhar Narayanan, and Christos Papadimitriou. 2024. On limitations of the transformer architecture. *arXiv preprint arXiv:2402.08164*.
- Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M. Sadler, Wei-Lun Chao, and Yu Su. 2023. [Llm-planner: Few-shot grounded planning for embodied agents with large language models](#).
- Kaya Stechly, Karthik Valmeekam, and Subbarao Kambhampati. 2024. [Chain of thoughtlessness? an analysis of cot in planning](#).
- Hrvoje Stojic, Eran Eldar, Hassan Bassam, Peter Dayan, and Raymond Dolan. 2018. Are you sure about that? on the origins of confidence in concept learning. In *Proceedings of the Cognitive Computational Neuroscience Conference*, pages 55–63.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Karthik Valmeekam, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati. 2022. Large language models still can't plan (a benchmark for llms on planning and reasoning about change). *arXiv preprint arXiv:2206.10498*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. [Chain-of-thought prompting elicits reasoning in large language models](#).
- Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze Lou, Yuandong Tian, Yanghua Xiao, and Yu Su. 2024. [Travelplanner: A benchmark for real-world planning with language agents](#). *arXiv preprint arXiv:2402.01622*.
- Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022a. [Webshop: Towards scalable real-world web interaction with grounded language agents](#). *Advances in Neural Information Processing Systems*, 35:20744–20757.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. [Tree of thoughts: Deliberate problem solving with large language models](#).
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2022b. [React: Synergizing reasoning and acting in language models](#). *arXiv preprint arXiv:2210.03629*.
- Roohfar Yousefzadeh and Jessica A. Mollick. 2021. [Extrapolation frameworks in cognitive psychology suitable for study of image classification models](#).
- Yanli Zhao, Andrew Gu, Rohan Varma, Liang Luo, Chien-Chin Huang, Min Xu, Less Wright, Hamid Shojanazeri, Myle Ott, Sam Shleifer, Alban Desmaison, Can Balioglu, Pritam Damania, Bernard Nguyen, Geeta Chauhan, Yuchen Hao, Ajit Mathews, and Shen Li. 2023. [Pytorch fsdp: Experiences on scaling fully sharded data parallel](#).
- Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. 2023. [Language agent tree search unifies reasoning acting and planning in language models](#). *arXiv preprint arXiv:2310.04406*.

A Experimental Details

A.1 Input detail

Table 4 describes a sample input of the model, describing the instruction of the Gridworld and the specific world information.

Common Prompt	<p>human: You are given a rectangular gridworld, where you can move up, down, left, or right as long as each of your x, y coordinates are within 0 to the x, y size of the grid. If you move up, your y coordinate increases by 1. If you move down, your y coordinate decreases by 1. If you move left, your x coordinate decreases by 1. If you move right, your x coordinate increases by 1.</p> <p>You will interact with the gridworld environment to reach the goal state, while avoiding the pit and the wall. You cannot move through the wall or move outside the grid. If you fall into the pit, you lose. If you reach the goal, you win. For each of your turn, you will be given the possible moves.</p> <p>You should respond your move with either one of 'up', 'down', 'left', or 'right'.</p> <p>gpt: OK</p> <p>human: Grid is from (0, 0) to (3, 2). Goal: (3, 2)\nCurrent: (0, 0)\nThe pit is at (3, 0). The wall is at (1, 1), and (2, 2).\nCurrent:\n(0, 0)\nPossible:\n(0, 1)\nup\n(1, 0)\nright</p>
---------------	--

Table 4: The prompt for all cases

A.2 Experimental setup details

We use one 8 Nvidia A100 node for both training and inference. For the training steps, we use FSDP framework (Zhao et al., 2023) and cosine annealing learning rate scheduler (Loshchilov and Hutter, 2017) for 1 epoch. We utilize bfloat16 floating-point format and a warmup ratio of 0.03. We set the weight decay as 0. We set the batch size of 2 for each GPUs, so the effective batch size is 16 per step. We train each model for $50000/16 = 3125$ steps. For inference, we use VLLM framework (Kwon et al., 2023).

While exploring the pure planning ability of the language model, we did not want the model to refuse to explore extrapolated data only because it has never seen the coordinate. To handle the bias, we adjust the starting position of the grid using uniform sampling while ensuring that the entire grid, with its x and y coordinates, fits within the range of 0 to 19. This method, illustrated in Figure 5, minimizes bias related to the unseen x and y coordinates by randomizing the starting point within the defined bounds.

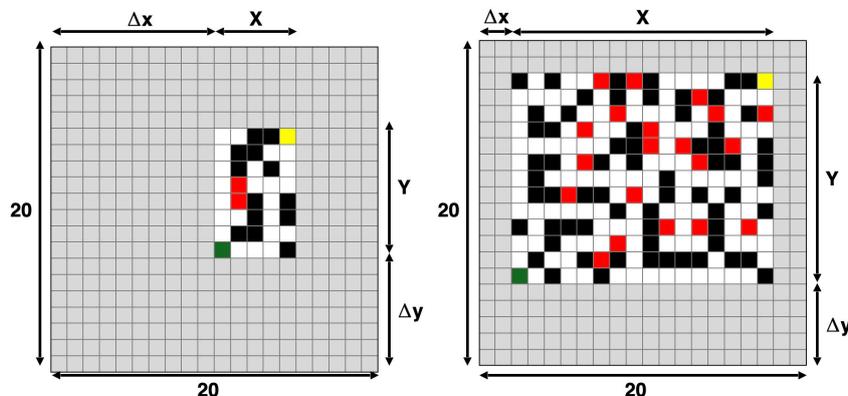


Figure 5: Visualization of configuring Gridworld instance for train(left) and test(right) dataset. To evaluate the extrapolation ability, we set the size of the grid as $X, Y \sim Unif(2, 10)$ for train and $X, Y \sim Unif(2, 20)$ for test. Also we set the starting point of the grid as $\Delta x \sim Unif(0, 20 - X)$, $\Delta y \sim Unif(0, 20 - Y)$ for both train and test.

B Cognitive map description

B.1 Cognitive map example: FWD

Table 5 describes a sample of FWD cognitive map construction for each experiment.

Design choice	Cognitive map example
NONE	
COT	Thought:\nStep 1:\nStep 2:\nStep 3:\nStep 4:\nStep 5:\nBacktrack:\n(3, 2)\nup\n(3, 1)\nrigh\n(2, 1)\nup\n(2, 0)\nrigh\n(1, 0)\nrigh\n(0, 0)
w.o. ALL BACKTRACK	Thought:\nStep 1:\n(0, 1)\nup\n(1, 0)\nrigh\nStep 2:\n(0, 2)\nup\n(2, 0)\nrigh\nStep 3:\n(1, 2)\nrigh\n(2, 1)\nup\nStep 4:\n(3, 1)\nrigh\nStep 5:\n(3, 2)\nup
w.o. ALL	Thought:\nStep 1:\n(0, 1)\nup\n(1, 0)\nrigh\nStep 2:\n(0, 2)\nup\n(2, 0)\nrigh\nStep 3:\n(1, 2)\nrigh\n(2, 1)\nup\nStep 4:\n(3, 1)\nrigh\nStep 5:\n(3, 2)\nup\nBacktrack:\n(3, 2)\nup\n(3, 1)\nrigh\n(2, 1)\nup\n(2, 0)\nrigh\n(1, 0)\nrigh\n(0, 0)
w.o. MARKING BACKTRACK	Thought:\nStep 1:\n(0, 1)\nup\n(0, -1)\ndown\n(-1, 0)\nleft\n(1, 0)\nrigh\nStep 2:\n(0, 2)\nup\n(0, 0)\ndown\n(-1, 1)\nleft\n(1, 1)\nrigh\n(1, 1)\nup\n(1, -1)\ndown\n(0, 0)\nleft\n(2, 0)\nrigh\nStep 3:\n(0, 3)\nup\n(0, 1)\ndown\n(-1, 2)\nleft\n(1, 2)\nrigh\n(2, 1)\nup\n(2, -1)\ndown\n(1, 0)\nleft\n(3, 0)\nrigh\nStep 4:\n(1, 3)\nup\n(1, 1)\ndown\n(0, 2)\nleft\n(2, 2)\nrigh\n(2, 2)\nup\n(2, 0)\ndown\n(1, 1)\nleft\n(3, 1)\nrigh\nStep 5:\n(3, 2)\nup\n(3, 0)\ndown\n(2, 1)\nleft\n(4, 1)\nrigh
w.o. MARKING	Thought:\nStep 1:\n(0, 1)\nup\n(0, -1)\ndown\n(-1, 0)\nleft\n(1, 0)\nrigh\nStep 2:\n(0, 2)\nup\n(0, 0)\ndown\n(-1, 1)\nleft\n(1, 1)\nrigh\n(1, 1)\nup\n(1, -1)\ndown\n(0, 0)\nleft\n(2, 0)\nrigh\nStep 3:\n(0, 3)\nup\n(0, 1)\ndown\n(-1, 2)\nleft\n(1, 2)\nrigh\n(2, 1)\nup\n(2, -1)\ndown\n(1, 0)\nleft\n(3, 0)\nrigh\nStep 4:\n(1, 3)\nup\n(1, 1)\ndown\n(0, 2)\nleft\n(2, 2)\nrigh\n(2, 2)\nup\n(2, 0)\ndown\n(1, 1)\nleft\n(3, 1)\nrigh\nStep 5:\n(3, 2)\nup\n(3, 0)\ndown\n(2, 1)\nleft\n(4, 1)\nrigh\nBacktrack:\n(3, 2)\nup\n(3, 1)\nrigh\n(2, 1)\nup\n(2, 0)\nrigh\n(1, 0)\nrigh\n(0, 0)
w.o. BACKTRACK	Thought:\nStep 1:\n(0, 1)\nup\n(0, -1)\ncut\n(-1, 0)\ncut\n(1, 0)\nrigh\nStep 2:\n(0, 2)\nup\n(0, 0)\ncut\n(-1, 1)\ncut\n(1, 1)\ncut\n(1, 1)\ncut\n(1, -1)\ncut\n(0, 0)\ncut\n(2, 0)\nrigh\nStep 3:\n(0, 3)\ncut\n(0, 1)\ncut\n(-1, 2)\ncut\n(1, 2)\nrigh\n(2, 1)\nup\n(2, -1)\ncut\n(1, 0)\ncut\n(3, 0)\ncut\nStep 4:\n(1, 3)\ncut\n(1, 1)\ncut\n(0, 2)\ncut\n(2, 2)\ncut\n(2, 2)\ncut\n(2, 0)\ncut\n(1, 1)\ncut\n(3, 1)\nrigh\nStep 5:\n(3, 2)\nup\n(3, 0)\ncut\n(2, 1)\ncut\n(4, 1)\ncut
MARKING deadend	Thought:\nStep 1:\n(0, 1)\nup\n(0, -1)\ncut\n(-1, 0)\ncut\n(1, 0)\nrigh\nStep 2:\n(0, 2)\nup\n(0, 0)\ncut\n(-1, 1)\ncut\n(1, 1)\ncut\n(1, 1)\ncut\n(1, -1)\ncut\n(0, 0)\ncut\n(2, 0)\nrigh\nStep 3:\n(0, 3)\ncut\n(0, 1)\ncut\n(-1, 2)\ncut\n(1, 2)\nrigh\n(2, 1)\nup\n(2, -1)\ncut\n(1, 0)\ncut\n(3, 0)\ncut\nStep 4:\n(1, 3)\ncut\n(1, 1)\ncut\n(0, 2)\ncut\n(2, 2)\ncut\n(2, 2)\ncut\n(2, 0)\ncut\n(1, 1)\ncut\n(3, 1)\nrigh\nStep 5:\n(3, 2)\nup\n(3, 0)\ncut\n(2, 1)\ncut\n(4, 1)\ncut\nBacktrack:\n(3, 2)\nup\n(3, 1)\nrigh\n(2, 1)\nup\n(2, 0)\nrigh\n(1, 0)\nrigh\n(0, 0)

Table 5: FWD cognitive map example for each experiment

B.2 Cognitive map example: BWD

Table 6 describes a sample of BWD cognitive map construction for each experiment.

Design choice	Cognitive map example
NONE	
COT	Thought:\nStep 1:\nStep 2:\nStep 3:\nStep 4:\nStep 5:\nBacktrack:\n(0, 0)\nright\n(1, 0)\nright\n(2, 0)\nup\n(2, 1)\nright\n(3, 1)\nup\n(3, 2)
w.o. ALL BACKTRACK	Thought:\nStep 1:\n(3, 1)\nup\nStep 2:\n(2, 1)\nright\nStep 3:\n(2, 0)\nup\nStep 4:\n(1, 0)\nright\nStep 5:\n(0, 0)\nright
w.o. ALL	Thought:\nStep 1:\n(3, 1)\nup\nStep 2:\n(2, 1)\nright\nStep 3:\n(2, 0)\nup\nStep 4:\n(1, 0)\nright\nStep 5:\n(0, 0)\nright\nBacktrack:\n(0, 0)\nright\n(1, 0)\nright\n(2, 0)\nup\n(2, 1)\nright\n(3, 1)\nup\n(3, 2)
w.o. MARKING BACKTRACK	Thought:\nStep 1:\n(3, 3)\ndown\n(3, 1)\nup\n(2, 2)\nright\n(4, 2)\nleft\nStep 2:\n(3, 2)\ndown\n(3, 0)\nup\n(2, 1)\nright\n(4, 1)\nleft\nStep 3:\n(2, 2)\ndown\n(2, 0)\nup\n(1, 1)\nright\n(3, 1)\nleft\nStep 4:\n(2, 1)\ndown\n(2, -1)\nup\n(1, 0)\nright\n(3, 0)\nleft\nStep 5:\n(1, 1)\ndown\n(1, -1)\nup\n(0, 0)\nright\n(2, 0)\nleft
w.o. MARKING	Thought:\nStep 1:\n(3, 3)\ndown\n(3, 1)\nup\n(2, 2)\nright\n(4, 2)\nleft\nStep 2:\n(3, 2)\ndown\n(3, 0)\nup\n(2, 1)\nright\n(4, 1)\nleft\nStep 3:\n(2, 2)\ndown\n(2, 0)\nup\n(1, 1)\nright\n(3, 1)\nleft\nStep 4:\n(2, 1)\ndown\n(2, -1)\nup\n(1, 0)\nright\n(3, 0)\nleft\nStep 5:\n(1, 1)\ndown\n(1, -1)\nup\n(0, 0)\nright\n(2, 0)\nleft\nBacktrack:\n(0, 0)\nright\n(1, 0)\nright\n(2, 0)\nup\n(2, 1)\nright\n(3, 1)\nup\n(3, 2)
w.o. BACKTRACK	Thought:\nStep 1:\n(3, 3)\ncut\n(3, 1)\nup\n(2, 2)\ncut\n(4, 2)\ncut\nStep 2:\n(3, 2)\ncut\n(3, 0)\ncut\n(2, 1)\nright\n(4, 1)\ncut\nStep 3:\n(2, 2)\ncut\n(2, 0)\nup\n(1, 1)\ncut\n(3, 1)\ncut\nStep 4:\n(2, 1)\ncut\n(2, -1)\ncut\n(1, 0)\nright\n(3, 0)\ncut\nStep 5:\n(1, 1)\ncut\n(1, -1)\ncut\n(0, 0)\nright\n(2, 0)\ncut
MARKING deadend	Thought:\nStep 1:\n(3, 3)\ncut\n(3, 1)\nup\n(2, 2)\ncut\n(4, 2)\ncut\nStep 2:\n(3, 2)\ncut\n(3, 0)\ncut\n(2, 1)\nright\n(4, 1)\ncut\nStep 3:\n(2, 2)\ncut\n(2, 0)\nup\n(1, 1)\ncut\n(3, 1)\ncut\nStep 4:\n(2, 1)\ncut\n(2, -1)\ncut\n(1, 0)\nright\n(3, 0)\ncut\nStep 5:\n(1, 1)\ncut\n(1, -1)\ncut\n(0, 0)\nright\n(2, 0)\ncut\nBacktrack:\n(0, 0)\nright\n(1, 0)\nright\n(2, 0)\nup\n(2, 1)\nright\n(3, 1)\nup\n(3, 2)

Table 6: BWD cognitive map example for each experiment

	w.o. ALL		with ALL	
Optimal	w.o. ALL BACKTRACK	w.o. ALL	w.o. MARKING	MARKING deadend
BWD	0.296	0.277	0.765	0.705
FWD	0.295	0.290	0.618	0.585
Reachable	w.o. ALL BACKTRACK	w.o. ALL	w.o. MARKING	MARKING deadend
BWD	0.394	0.283	0.724	0.885
FWD	0.416	0.345	0.816	0.854

Table 7: Planning performance with ALL exclusion

	w.o. BACKTRACK		with BACKTRACK	
Optimal	w.o. MARKING BACKTRACK	w.o. BACKTRACK	w.o. MARKING	MARKING deadend
BWD	0.406	0.423	0.765	0.705
FWD	0.528	0.516	0.618	0.585
Reachable	w.o. MARKING BACKTRACK	w.o. BACKTRACK	w.o. MARKING	MARKING deadend
BWD	0.739	0.852	0.724	0.885
FWD	0.672	0.624	0.816	0.854

Table 8: Planning performance with BACKTRACK exclusion

C Additional results

C.1 Additional investigations

Effect of ALL inclusion Both cognitive map constructions w.o. ALL (w.o. ALL BACKTRACK: 29.6% and 29.5% for BWD and FWD construction, respectively; w.o. ALL: 27.7% and 29.0% for BWD and FWD construction, respectively) suffer at generating the optimal plan, achieving under 30% success rate. However, every cognitive map construction with the inclusion of ALL shows better performance, with FWD MARKING deadend being the lowest among them(58.5%).

We observe a similar trend in the reachable planning test. While both baseline methods w.o. ALL performed at best 41.6% (FWD w.o. ALL BACKTRACK), the lowest performance among every cognitive map construction was 72.4% (BWD w.o. MARKING). This implies that the **inclusion of ALL significantly enhances the planning capability, leading to more successful and efficient pathfinding.**

Effect of BACKTRACK inclusion Both cognitive map constructions w.o. BACKTRACK (w.o. MARKING BACKTRACK: 40.6% and 52.8% for BWD and FWD construction, respectively; w.o. BACKTRACK: 42.3% and 51.6% for BWD and FWD construction, respectively) suffer at generating the optimal plan. However, every cognitive map construction with the inclusion of BACKTRACK shows better performance, with FWD MARKING deadend being the lowest among them(58.5%).

The analysis in the reachable planning test was slightly blurry, yet there was an obvious trend. For each experiment, adding backtracking enhanced the performance of the planning in most settings(except BWD construction w.o. MARKING). This implies that the **inclusion of BACKTRACK slightly enhances the planning capability.**

C.2 Vizualization for optimal planning experiments

For optimal planning, we have only success or failure cases. Hence we only provide the success rate for each experiment.

FWD construction See Figure 6 for success rate of each experiment.

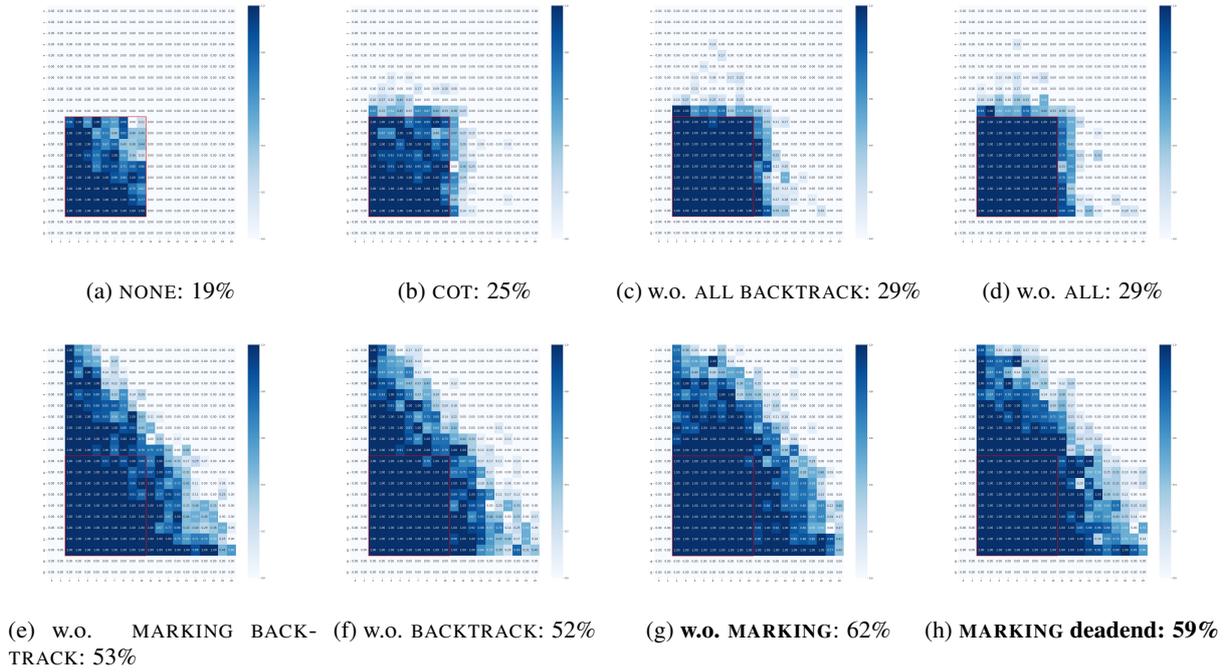


Figure 6: Success rate for optimal planning, FWD construction

BWD construction See Figure 7 for success rate of each experiment.



Figure 7: Success rate for optimal planning, BWD construction

C.3 Visualization for reachable planning experiments

For reachable planning, we have one success cases and three different failure cases(deadend, max step, and invalid). Hence we provide the visualization of all 4 cases.

FWD construction For success rate, see Figure 8. For failure cases, see Figure 9 for deadend, Figure 10 for max step, and Figure 11 for invalid rate.

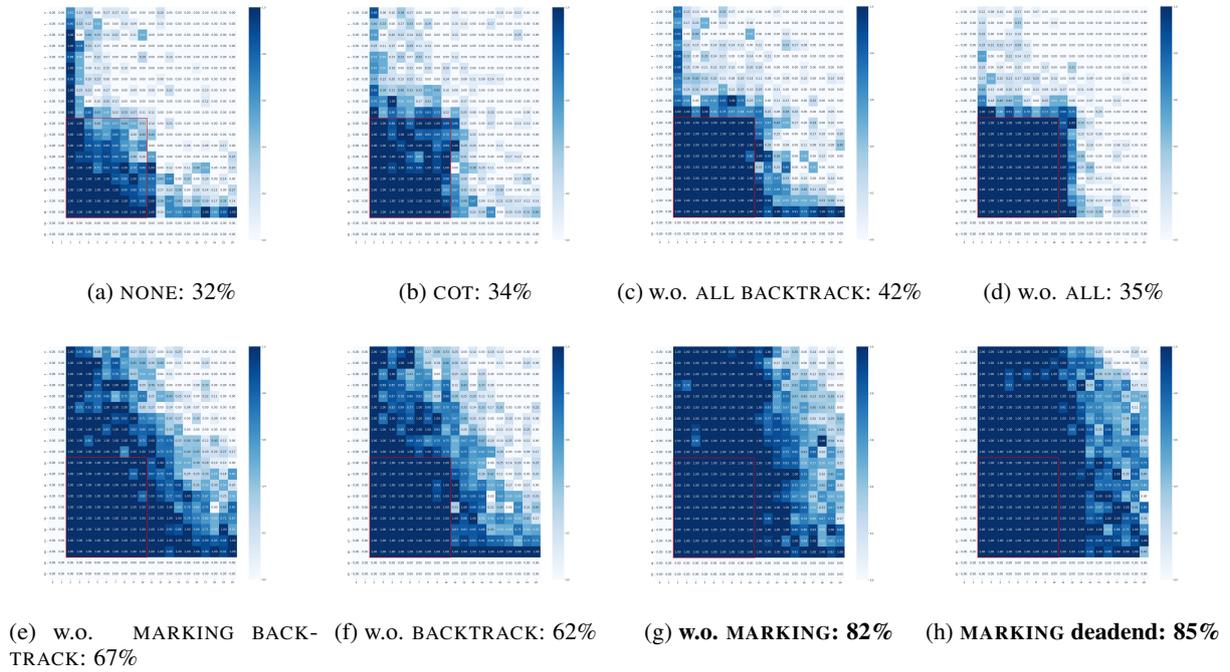


Figure 8: Success rate for reachable planning, FWD construction

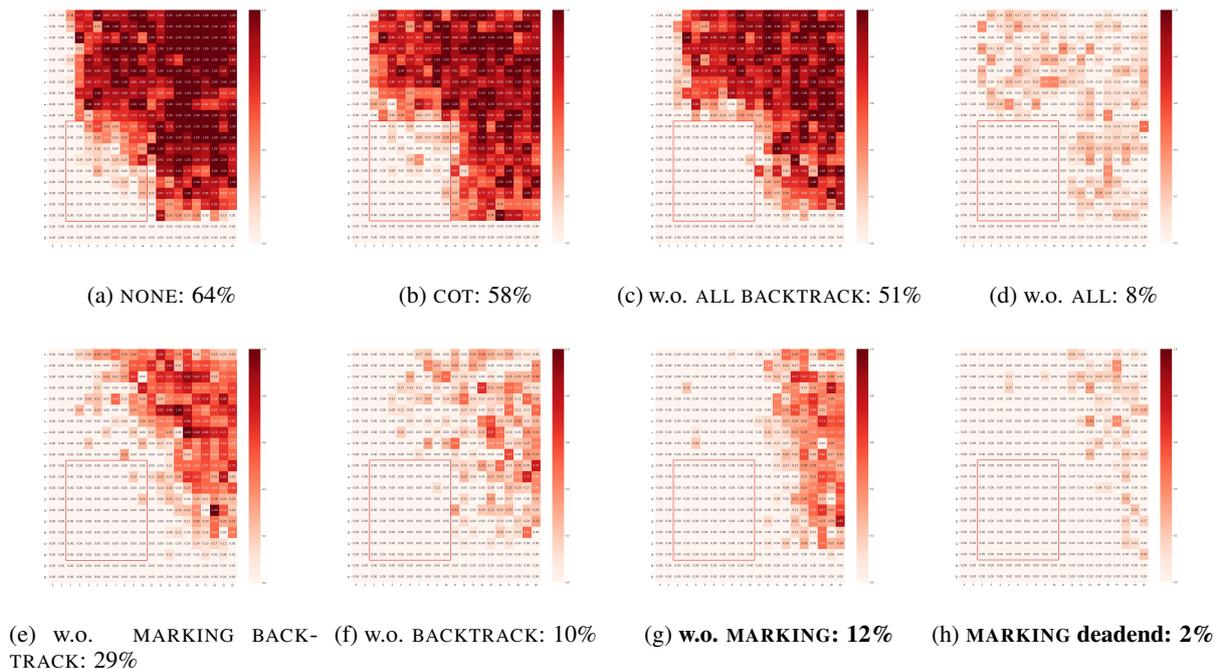
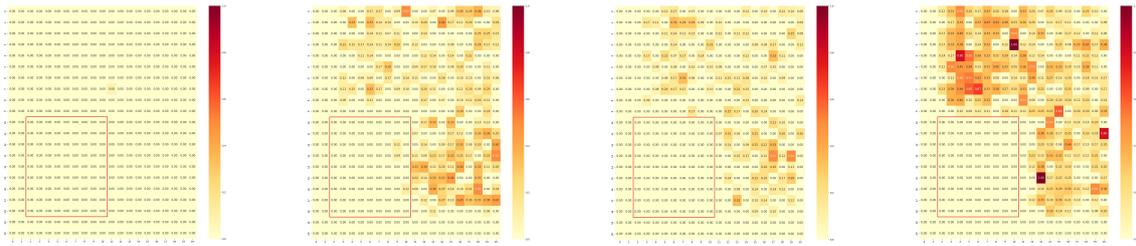


Figure 9: Deadend rate for reachable planning, FWD construction

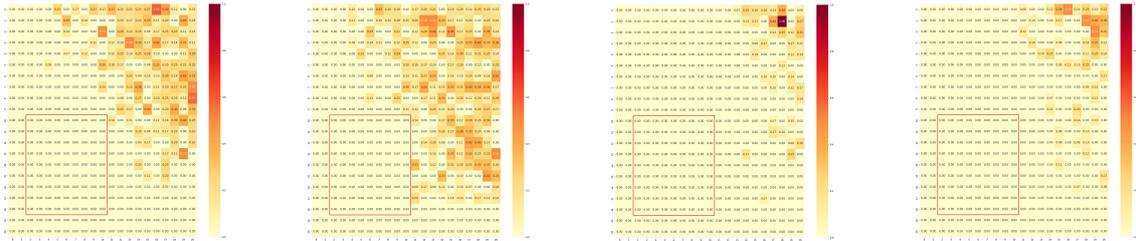


(a) NONE: 0%

(b) COT: 8%

(c) w.o. ALL BACKTRACK: 5%

(d) w.o. ALL: 15%



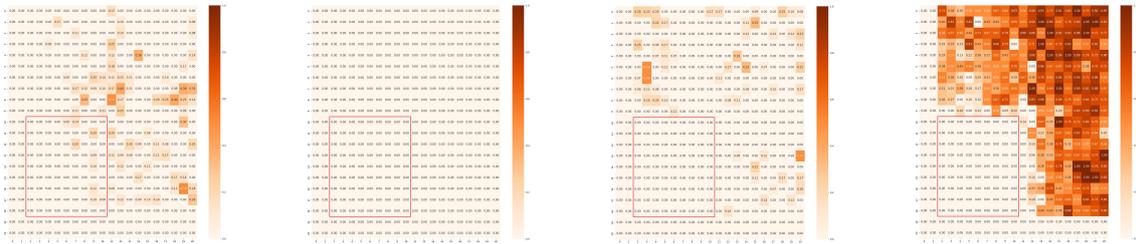
(e) w.o. MARKING TRACK: 7%

(f) w.o. BACKTRACK: 10%

(g) w.o. MARKING: 2%

(h) MARKING deadend: 3%

Figure 10: Max step rate for reachable planning, FWD construction

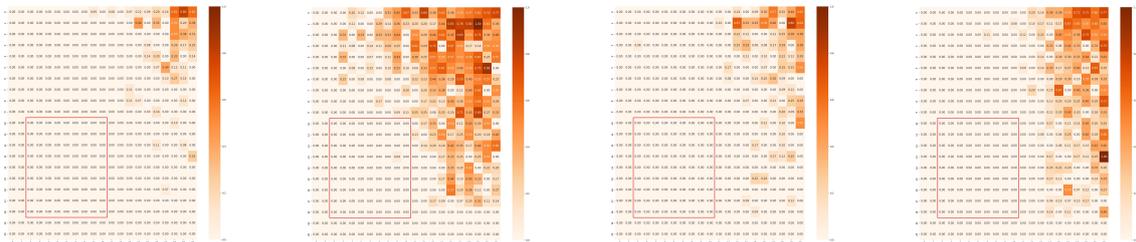


(a) NONE: 3%

(b) COT: 0%

(c) w.o. ALL BACKTRACK: 3%

(d) w.o. ALL: 43%



(e) w.o. MARKING TRACK: 3%

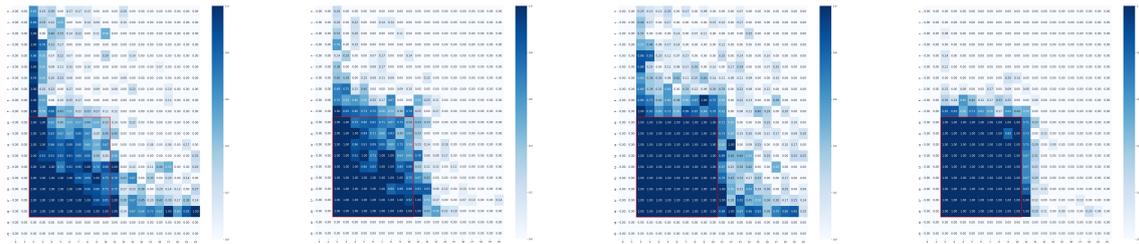
(f) w.o. BACKTRACK: 18%

(g) w.o. MARKING: 4%

(h) MARKING deadend: 10%

Figure 11: Invalid rate for reachable planning, FWD construction

BWD construction For success rate, see Figure 12. For failure cases, see Figure 13 for deadend, Figure 14 for max step, and Figure 15 for invalid rate.

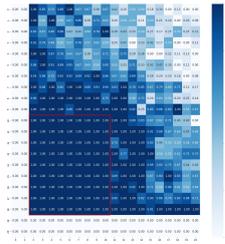


(a) NONE: 32%

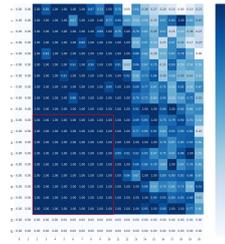
(b) COT: 29%

(c) w.o. ALL BACKTRACK: 40%

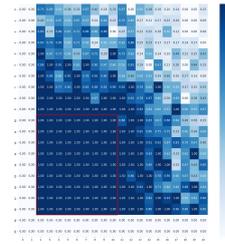
(d) w.o. ALL: 28%



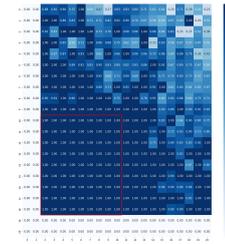
(e) w.o. MARKING TRACK: 74%



(f) w.o. BACKTRACK: 85%

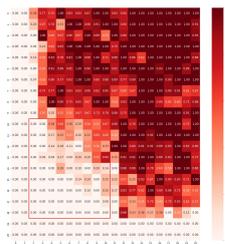


(g) w.o. MARKING: 72%

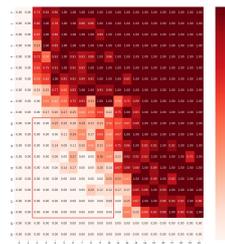


(h) MARKING deadend: 89%

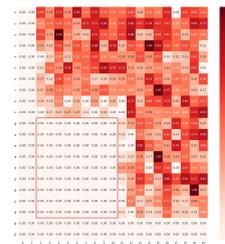
Figure 12: Success rate for reachable planning, BWD construction



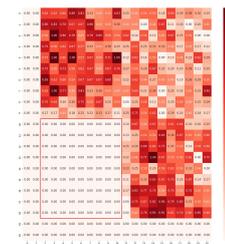
(a) NONE: 64%



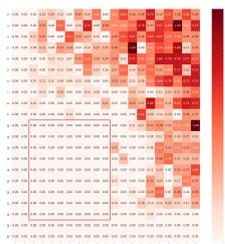
(b) COT: 70%



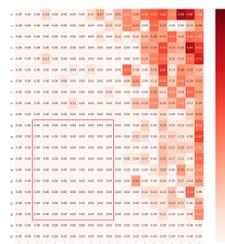
(c) w.o. ALL BACKTRACK: 33%



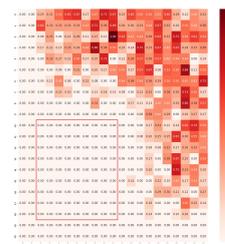
(d) w.o. ALL: 37%



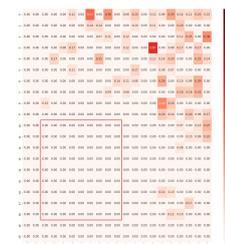
(e) w.o. MARKING TRACK: 19%



(f) w.o. BACKTRACK: 10%

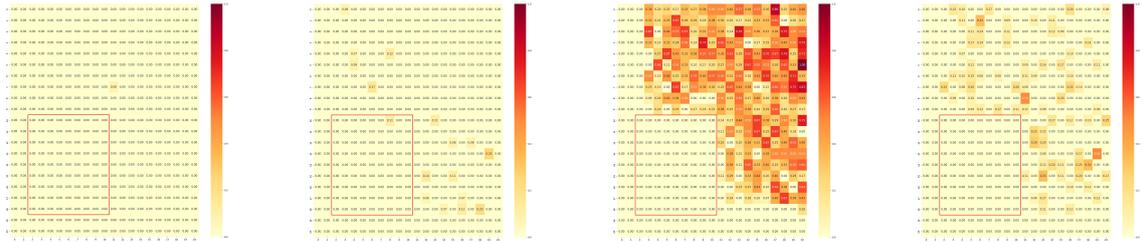


(g) w.o. MARKING: 20%



(h) MARKING deadend: 3%

Figure 13: Deadend rate for reachable planning, BWD construction

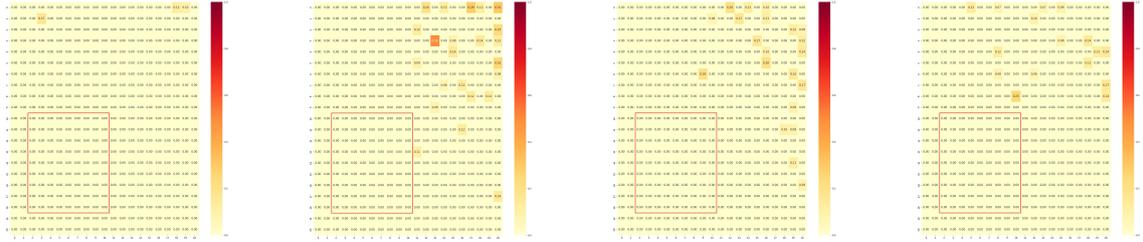


(a) NONE: 0%

(b) COT: 1%

(c) w.o. ALL BACKTRACK: 26%

(d) w.o. ALL: 3%



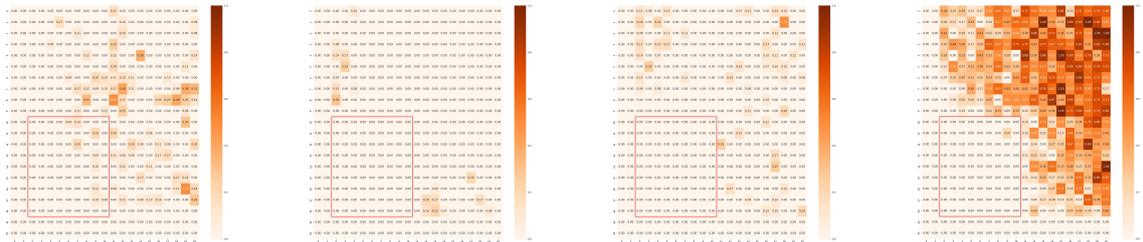
(e) w.o. MARKING BACKTRACK: 0%

(f) w.o. BACKTRACK: 1%

(g) w.o. MARKING: 1%

(h) MARKING deadend: 1%

Figure 14: Max step rate for reachable planning, BWD construction

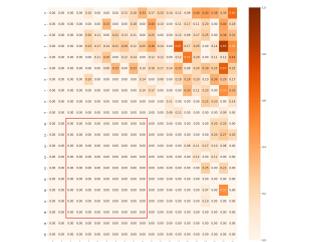


(a) NONE: 3%

(b) COT: 1%

(c) w.o. ALL BACKTRACK: 2%

(d) w.o. ALL: 32%



(e) w.o. MARKING BACKTRACK: 7%

(f) w.o. BACKTRACK: 4%

(g) w.o. MARKING: 7%

(h) MARKING deadend: 8%

Figure 15: Invalid rate for reachable planning, BWD construction

D Comparison of cognitive model with exploration-based model

See Figure 16 for comparing the success rate of reachable planning between the cognitive model and exploration-based planning. See Algorithm 2 for the algorithm of DFS. Moreover, Table 9 describes an example of an output of exploration-based planning.

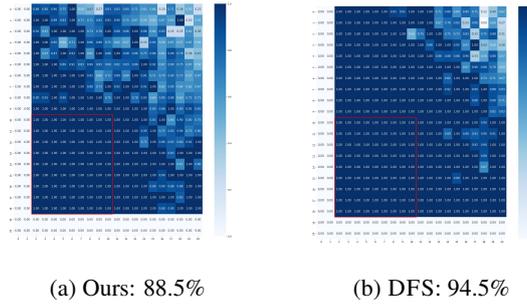


Figure 16: Reachability plan performance between ours(BWD MARKING deadend) and DFS

DFS Prompt

gpt: up
human: Current:\n(0, 1)\nPossible:\n(0, 2)\nup\n(0, 0)\ndown
gpt: up
human: Current:\n(0, 2)\nPossible:\n(0, 1)\ndown\n(1, 2)\nrigh
gpt: right
human: Current:\n(1, 2)\nPossible:\n(0, 2)\nleft
gpt: left
human: Current:\n(0, 2)\nPossible:\n(0, 1)\ndown\n(1, 2)\nrigh
gpt: down
human: Current:\n(0, 1)\nPossible:\n(0, 2)\nup\n(0, 0)\ndown
gpt: down
human: Current:\n(0, 0)\nPossible:\n(0, 1)\nup\n(1, 0)\nrigh
gpt: right
human: Current:\n(1, 0)\nPossible:\n(0, 0)\nleft\n(2, 0)\nrigh
gpt: right
human: Current:\n(2, 0)\nPossible:\n(2, 1)\nup\n(1, 0)\nleft
gpt: up
human: Current:\n(2, 1)\nPossible:\n(2, 0)\ndown\n(3, 1)\nrigh
gpt: right
human: Current:\n(3, 1)\nPossible:\n(3, 2)\nup\n(2, 1)\nleft
gpt: up

Table 9: The prompt for DFS case

Algorithm 2 Tree search via DFS

```

s ← start, queue ← [s]
while s ≠ goal do
  if |EXPLORE(s)| ≥ 1 then
    tmp ← RANDOM(EXPLORE(s))
    PARENT(tmp) ← s
    s ← tmp
  end if
  s ← PARENT(s)
  queue += s
end while

```

▷ If there is a sample from s not visited yet
