

---

# Can language model plan in extrapolated environments?: Casestudy in textualized Gridworld

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 While language models have demonstrated impressive capabilities across general-  
2 ized language tasks, their ability to extrapolate in a certain task is highly unknown.  
3 We first introduce the optimal path planning task in a textualized Gridworld environ-  
4 ment as a valid probe for estimating the extrapolability of language models. We  
5 show that the mere next token prediction inherently fails to extrapolate in solving  
6 the task. Inspired by human cognition, we claim that language models should  
7 construct an internal simulation that explores the environment, i.e. cognitive map  
8 before actually interacting with the given environment. We demonstrate that auto-  
9 regressive generation of cognitive map and planning sequence can significantly  
10 enhance the performance of the planning power even in extrapolated environments,  
11 suggesting the necessity of cognitive map for language models as a path forward.

## 12 1 Introduction

13 Language models have recently demonstrated remarkable proficiency in a variety of complex  
14 tasks (Brown et al., 2020; Touvron et al., 2023; Chowdhery et al., 2023; Chen et al., 2021), from  
15 natural language understanding to code generation, primarily through the training objective of next  
16 token prediction. This training paradigm enables language models to excel at general planning tasks  
17 by leveraging their extensive learned knowledge and pattern recognition capabilities (Ahn et al.,  
18 2022; Liang et al., 2023; Song et al., 2023). However, language model also falter in scenarios that  
19 require robust, long-horizon planning tasks (Dziri et al., 2024). It is in contrast that humans naturally  
20 employ model-based planning, internally construct models to simulate outcomes and guide optimal  
21 decision-making, as extensively documented in cognitive science literature (Daw et al., 2005). This  
22 distinction aligns with the dual-process theory of reasoning: System 1 processes are fast, automatic,  
23 and pattern-based, akin to model-free planning, while System 2 processes are slower, deliberative,  
24 and involve explicit reasoning (Daniel, 2017).

25 The Chain of Thought (CoT) approach (Wei et al., 2023) is a prominent method used by language  
26 models to emulate the System 2 cognitive process, optimizing the reasoning pathway. While CoT  
27 has significantly improved reasoning and planning, existing methods focus solely on generating  
28 intermediate steps from the initial state to the goal (Wei et al., 2023; Nye et al., 2021). This approach  
29 deviates from the concept of "simulation," which involves analyzing real-world systems and predicting  
30 outcomes. In contrast, human System 2 cognition typically engages in iterative simulations, refining  
31 steps until reaching the goal while avoiding dead-end states. After reaching the goal, humans often  
32 work backward to determine the optimal actions leading to that state. This process, which forms a  
33 "cognitive map," is controlled by the prefrontal cortex (Daw et al., 2005), allowing for deliberate,  
34 complex decision-making (Doll et al., 2012).

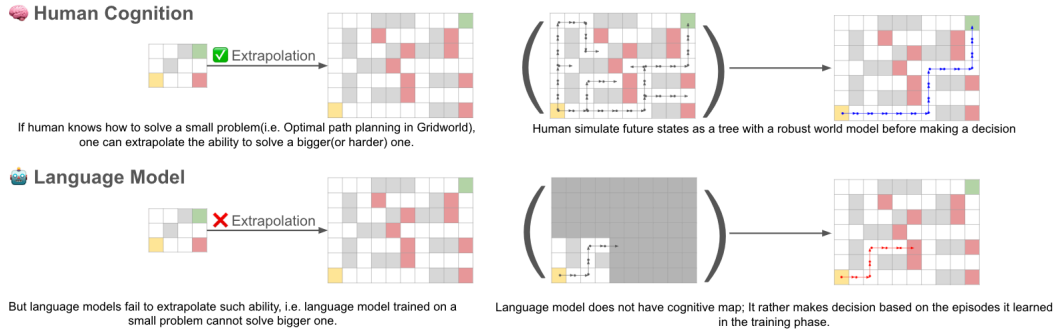


Figure 1: Difference between human cognition and language models: When both humans and language models are taught to path plan in small Gridworld environments, humans can successfully extrapolate to larger environments after a few demonstrations. This ability is due to the construction of a mental map of future states, known as a cognitive map. In contrast, language models struggle with extrapolation, primarily because they rely solely on the training data. We propose that language models should develop the capability to internally simulate future states, organized as a decision tree, to improve performance on tasks requiring extrapolation.

35 We argue that the lack of a cognitive map in current language models is a primary reason they fail to  
 36 adapt learned behaviors from training into extrapolated environments, as depicted in Figure 1. To  
 37 support this claim, we introduce an optimal path-planning task in a textualized Gridworld environment  
 38 as a benchmark for assessing language model extrapolability. Our results demonstrate that language  
 39 models struggle to extrapolate in such environments when trained purely through imitation learning  
 40 and even with CoT fine-tuning. As a solution, we insist that language models also need human-like  
 41 cognitive maps. To support the idea, we propose training language models using datasets augmented  
 42 with a simple cognitive map(See Figure 2). We find that auto-regressive generation of cognitive maps  
 43 and planning sequences can significantly improve the model’s ability to plan effectively, even in  
 44 unfamiliar, extrapolated environments. These results indicate that integrating cognitive maps into  
 45 language models may be a promising step toward achieving human-like cognition and enhancing  
 46 their capacity to generalize to more composite and unseen environments.

## 47 2 Experimental setup

### 48 2.1 Basic setup

49 In this paper, we set textualized Gridworld (Brown, 2015) as the main task. Gridworld is a task  
 50 that involves path planning from the start state to the goal state while avoiding the pit, wall, and  
 51 grids outside the world. Especially, we ensure that there is only one path from start to goal for each  
 52 environment. We prompt a textualized instruction of the Gridworld to the model and set a textworld  
 53 environment that receives an action(either up, down, left, or right) and outputs the corresponding  
 54 transition state(See Appendix B.1 for textualized input instruction).

55 We choose Gridworld because of following reasons: First, Gridworld requires minimal knowl-  
 56 edge. We can probe the extrapolability of language models with board game such as Einstein’s  
 57 puzzle (Brainzilla, 2017) or Blocksworld (Valmeekam et al., 2022), but with assessing minimal  
 58 world knowledge in order to probe the pure extrapolability of the model. Gridworld only requires  
 59 4 actions(up, down, right and left) and the transition of each action is simple and explicit, which  
 60 is a good fit. Second, we can generate Gridworld environment of an arbitrary size. After training  
 61 with few demonstrations, we need to find a larger environment which was unseen during the training  
 62 phase in order to test the extrapolability of the model. Gridworld, unlike other planning tasks such  
 63 as coding or math, can always generate a board of a bigger size that was unseen during the training  
 64 phase. This makes train and validation much easier. Last, Gridworld is outside  $TC^0$  complexity  
 65 class, making it impossible to solve with mere next token prediction. Existing works implies that  
 66 the task outside  $TC^0$  class cannot be solved directly with next token prediction with fixed precision  
 67 transformer architecture (Merrill and Sabharwal, 2023).

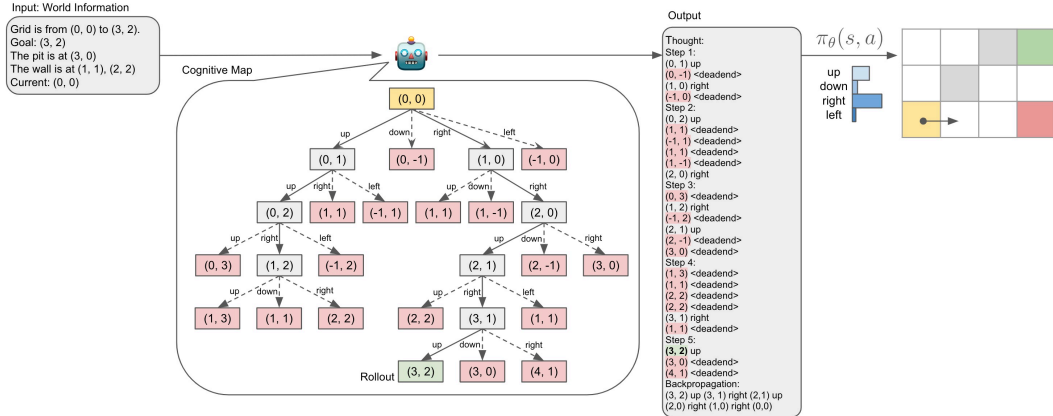


Figure 2: Optimal path-planning with cognitive map: 1. We initialize the world environment which the model is going to interact. 2. The textualized input instruction containing information about the environment( $u$ ) is fed into the model. 3. Before interacting with the world, the model constructs a textualized cognitive map( $m$ ). Our construction of  $m$  is a tree-structured verbal representation of the world model, which is represented in a sequential manner. 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. We show that the cognitive map deduces optimal plan, and it shows human-cognitive characteristics such as generalization to extrapolated environments(Section 3).

68 For each observation, we give the information of the current state along with possible moves that do  
69 not directly result the deadend state along with its corresponding states. For example, if the current  
70 state  $s_t$  is (11, 4) and there the possible actions are right or left, the observation  $o_t$  is "Current:\n(11,  
71 4)\nPossible:\n(10, 4)\nleft\n(12, 4)\nrigh".(See Appendix B.2 for details)

## 72 2.2 Baselines

73 We name NONE and COT as our baseline experiments representing imitation-based planning. NONE  
74 implicitly learns how to conduct path planning, while COT learns a verbalized backward trace  
75 as a Chain of Thought(CoT) demonstration. For example, the COT reasoning of the Figure 1 is  
76 "(3, 2)up\n(3, 1)\nrigh\n(2, 1)\nup\n(2, 0)\nrigh\n(1, 0)\nrigh\n(0, 0)", while we do not verbalize  
77 anything for NONE. See Appendix E.1 for the whole construction.

## 78 2.3 Cognitive map design choices

79 The design of the cognitive map involves three key processes: **Sampling**, **Propagation**, and **Back-**  
80 **tracking**. The model begins by sampling plausible actions at each state and then propagates forward  
81 by exploring new states that result from these actions. This process of sampling and propagating  
82 continues until the model reaches the goal. Once the goal is achieved, the model backtracks from the  
83 goal state to the starting state to refine the optimal path(See Figure 2).

84 Our construction of the cognitive map is straightforward; We sample all 4 possible directions("up",  
85 "down", "right", and "left") as actions, and we propagate each action iteratively starting from the goal  
86 state until reaching the start state. After reaching the start state, we backtrack until reaching the goal  
87 state. We comprise these steps as Sampling, Propagation and Backtracking stage. See Appendix C  
88 for more details. We See Appendix D for ablations of different cognitive map constructions.

## 89 2.4 Experiment scheme

90 We observe the capability of the language model when generating the optimal plan without any further  
91 observations for each experiments, NONE, COT and COGNITIVE MAP. We evaluate the optimality  
92 of the trajectory generated by  $a_1, a_2, \dots, a_n \sim \pi_{\theta}(\cdot|u, m) \in \mathcal{A}^n$ (For preliminary notation, see  
93 Appendix A). The only difference among the experiments is the generation of  $m$ . For NONE, we do  
94 not generate  $m$ . For COT, we generate a backward trace from the goal to the start. For COGNITIVE  
95 MAP, we generate a simulation process from goal to the start.

96 Given the instruction  $u$  and generated (action, state) trajectory  $k$ , we first define optimal path  $k^*$  as  
 97 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  
 98 need to check whether the generated trajectory is  $k^*$ . Now we can define  $R(u, k)$  as follows:

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

101 We compute the mean value of  $R(u, k)$  across each size of the Gridworld to probe the extrapolability  
 102 of the language model. If we get a high value beyond the training boundary, that means that the  
 103 model has a good extrapolability, otherwise it cannot extrapolate into bigger environments.

### 104 3 Result and Analysis

#### 105 3.1 Cognitive map induces extrapolability

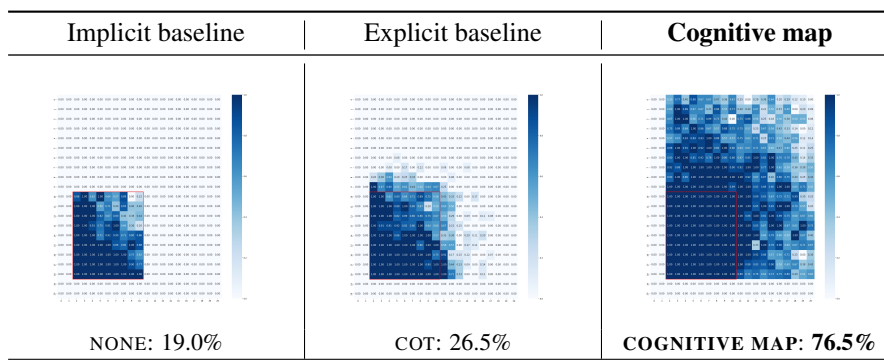


Table 1: Qualitative comparison of optimal plan generation rate between baselines and our methods. The degree of the darkness at  $(x, y)$  coordinate of each plot denotes the performance of the corresponding model in Gridworld of size  $x \times y$ . The red box denotes the boundary of the training data. We provide visualization of all the experiments in Appendix F.

106 As stated in Section 2.1, we train the model on world sizes of up to  $10 \times 10$  and test it on world sizes  
 107 of up to  $20 \times 20$  world to investigate the extrapolation ability. Extrapolation success is defined as the  
 108 ability to succeed in rollouts beyond the  $10 \times 10$  boundary. As shown in Table 1, the darkness outside  
 109 the boundary of the red box shows that **the cognitive map helps in planning on the extrapolated**  
 110 **data**. We also observe a consistent tendency for experiments with ALL inclusion, which have a wider  
 111 coverage of success rate in the extrapolated data. These results align with the result discussed in  
 112 Section 3.1. We analyze ablation results of cognitive map in Appendix F.

#### 113 3.2 Why does cognitive map induce extrapolation?

114 A language model can be viewed as a statistical model trained to minimize the KL divergence between  
 115 its predicted logits and the training objective. Given the training data distribution  $P_{train}$ , the model  
 116 is optimized to fit this distribution, enabling it to predict interpolated data points. However, the  
 117 model’s ability to extrapolate beyond the training data is less understood and lacks a clear theoretical  
 118 explanation. So it is not surprising for NONE and COT not having extrapolability. But what is the key  
 119 difference of cognitive map that enables such ability? Our current insight is as follows:

120 Human demonstrations are often abstract and entangled, with the human policy  $\pi_{\text{human}}$  being inter-  
 121 twined with the observed replication policy  $\pi_{\text{human demo}}$ . The replication policy tends to be more  
 122 complex due to hidden variables and factors influencing human demonstrations. We hypothesize  
 123 that the human policy  $\pi_{\text{human}}$  is relatively simple and can be modeled using a combination of three  
 124 key functions: sampling ( $S$ ), propagation ( $T$ ), and backtracking ( $T^{-1}$ ). By leveraging these simple,  
 125 well-defined functions, the model can effectively handle path-planning tasks and demonstrate robust  
 126 extrapolative abilities beyond the training data.



127 **References**

- 128 Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea  
129 Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Daniel Ho, Jasmine  
130 Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Eric Jang, Rosario Jauregui Ruano, Kyle Jeffrey, Sally  
131 Jesmonth, Nikhil J Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Kuang-Huei Lee,  
132 Sergey Levine, Yao Lu, Linda Luu, Carolina Parada, Peter Pastor, Jornell Quiambao, Kanishka  
133 Rao, Jarek Rettinghouse, Diego Reyes, Pierre Sermanet, Nicolas Sievers, Clayton Tan, Alexander  
134 Toshev, Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Mengyuan Yan, and Andy  
135 Zeng. 2022. Do as i can, not as i say: Grounding language in robotic affordances.
- 136 Brainzilla. 2017. Einstein’s riddle. Accessed on September 10, 2024.
- 137 Brandon Brown. 2015. Q-learning with neural networks: Learning gridworld with q-learning.  
138 Accessed on June 19, 2024.
- 139 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
140 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models  
141 are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- 142 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared  
143 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating  
144 large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- 145 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam  
146 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm:  
147 Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.
- 148 Kahneman Daniel. 2017. *Thinking, fast and slow*.
- 149 Nathaniel D Daw, Yael Niv, and Peter Dayan. 2005. Uncertainty-based competition between prefrontal  
150 and dorsolateral striatal systems for behavioral control. *Nature neuroscience*, 8(12):1704–1711.
- 151 Bradley B Doll, Dylan A Simon, and Nathaniel D Daw. 2012. The ubiquity of model-based  
152 reinforcement learning. *Current opinion in neurobiology*, 22(6):1075–1081.
- 153 Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jiang, Bill Yuchen Lin, Sean  
154 Welleck, Peter West, Chandra Bhagavatula, Ronan Le Bras, et al. 2024. Faith and fate: Limits of  
155 transformers on compositionality. *Advances in Neural Information Processing Systems*, 36.
- 156 Mehran Kazemi, Najoung Kim, Deepti Bhatia, Xin Xu, and Deepak Ramachandran. 2023. Lambada:  
157 Backward chaining for automated reasoning in natural language.
- 158 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.  
159 Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language  
160 model serving with pagedattention.
- 161 Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and  
162 Andy Zeng. 2023. Code as policies: Language model programs for embodied control.
- 163 Ilya Loshchilov and Frank Hutter. 2017. Sgdr: Stochastic gradient descent with warm restarts.
- 164 William Merrill and Ashish Sabharwal. 2023. The parallelism tradeoff: Limitations of log-precision  
165 transformers.
- 166 Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David  
167 Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, et al. 2021. Show your work:  
168 Scratchpads for intermediate computation with language models. *arXiv preprint arXiv:2112.00114*.
- 169 Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M. Sadler, Wei-Lun Chao, and Yu Su. 2023.  
170 Llm-planner: Few-shot grounded planning for embodied agents with large language models.
- 171 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
172 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open  
173 and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

174 Karthik Valmeekam, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati. 2022. Large  
175 language models still can't plan (a benchmark for llms on planning and reasoning about change).  
176 [arXiv preprint arXiv:2206.10498](#).

177 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le,  
178 and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models.

179 Yanli Zhao, Andrew Gu, Rohan Varma, Liang Luo, Chien-Chin Huang, Min Xu, Less Wright,  
180 Hamid Shojanazeri, Myle Ott, Sam Shleifer, Alban Desmaison, Can Balioglu, Pritam Damania,  
181 Bernard Nguyen, Geeta Chauhan, Yuchen Hao, Ajit Mathews, and Shen Li. 2023. Pytorch fsdp:  
182 Experiences on scaling fully sharded data parallel.

## 183 A Planning task interacting with world model

184 Planning task with environment feedback in language model can be formalized as a Markov decision  
185 process with instruction space  $\mathcal{U}$ , state space  $\mathcal{S}$ , action space  $\mathcal{A}$ , observation space  $\mathcal{O}$ , metric space  
186  $\mathcal{C}$ , 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  
187 the trajectory length. Note that for language model domain,  $\mathcal{U}$ ,  $\mathcal{A}$ , and  $\mathcal{O}$  are given as sequences of  
188 language tokens.

189 Given an instruction  $u \in \mathcal{U}$ , the model  $\theta$  first generates the action  $a_1 \sim \pi_\theta(\cdot|u) \in \mathcal{A}$  according to its  
190 policy  $\pi_\theta$ . For each state  $s_t \in \mathcal{S}$  and its observation  $o_t \in \mathcal{O}$ , the agent generates the corresponding  
191 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  
192  $s_{t+1} = T(s_t, a_t)$  and its observation  $o_{t+1}$ . The interaction loop repeats until the task has terminated  
193 for some reason(succeeded, failed, number of steps exceeded maximum value, etc.), and the action  
194 trajectory is denoted as:

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

195

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

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

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

## 201 B Experimental Details

### 202 B.1 Input detail

203 Table 2 describes a sample input of the model, describing the instruction of the Gridworld and the  
204 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.\n\nYou 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.\n\nYou 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 2: The prompt for all cases

205 **B.2 Experimental setup details**

206 We train the model for 1 epoch with 50K samples of size  $10 \times 10$  at largest. After training, we test  
 207 the model with 3K samples with each grid being  $20 \times 20$  at largest. We utilize Llama-3-8B model<sup>1</sup>  
 208 throughout the whole experiments. In each turn, we set the maximum token length of the model to be  
 209 8192.

210 We use one 8 Nvidia A100 node for both training and inference. For the training steps, we use FSDP  
 211 framework (Zhao et al., 2023) and cosine annealing learning rate scheduler (Loshchilov and Hutter,  
 212 2017) for 1 epoch. We utilize bfloat16 floating-point format and a warmup ratio of 0.03. We set the  
 213 weight decay as 0. We set the batch size of 2 for each GPUs, so the effective batch size is 16 per step.  
 214 We train each model for  $50000/16 = 3125$  steps. For inference, we use VLLM framework (Kwon  
 215 et al., 2023).

216 While exploring the pure planning ability of the language model, we did not want the model to refuse  
 217 to explore extrapolated data only because it has never seen the coordinate. To handle the bias, we  
 218 adjust the starting position of the grid using uniform sampling while ensuring that the entire grid,  
 219 with its x and y coordinates, fits within the range of 0 to 19. This method, illustrated in Figure 3,  
 220 minimizes bias related to the unseen x and y coordinates by randomizing the starting point within the  
 221 defined bounds.

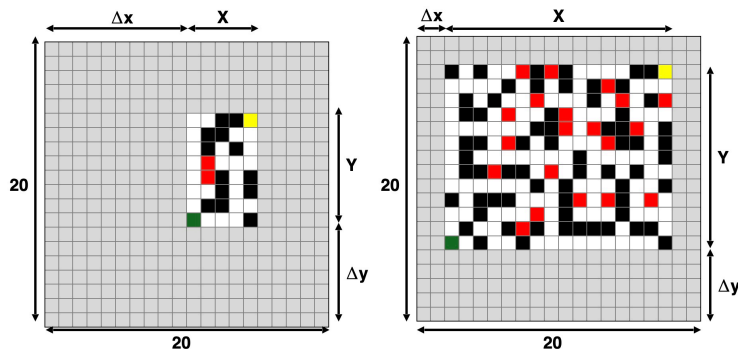


Figure 3: Visualization of configuring Gridworld instance for train(left) and test(right) dataset. To evaluate the extrapolation ability, we set the size of the grid as  $X, Y \sim Unif(2, 10)$  for train and  $X, Y \sim Unif(2, 20)$  for test. Also we set the starting point of the grid as  $\Delta x \sim Unif(0, 20 - X)$ ,  $\Delta y \sim Unif(0, 20 - Y)$  for both train and test.

<sup>1</sup><https://llama.meta.com/llama3/>

## 222 C Cognitive map for language model

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

### 228 C.1 Sampling

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

### 233 C.2 Propagation

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

### 240 C.3 Backtrack

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

### 249 C.4 Training how to generate cognitive map

250 To sum up, constructing the cognitive map  $m$  is a sequential application of sampling( $S$ ),  
251 propagation( $T$ ), and backtrack( $T^{-1}$ ). We train the language model  $\theta$  via supervised learning method  
252 so that the model can successfully construct the cognitive map without any external interaction, as  
253 shown in Figure 2.

## 254 D Experiment Ablation

### 255 D.1 Reachable planning analysis: multi-turn setting

256 We also investigate whether the cognitive map can deduce a reachable plan(which doesn’t require  
257 optimal planning) in a multi-turn setting. We evaluate the reachability of the trajectory  $a_1, a_2, \dots, a_n$   
258 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  
259 the model and the environment in a multi-turn setting, so failure cases can be divided into three  
260 cases(deadend, max step, and invalid). Especially, given the instruction  $u$  and generated (action, state)  
261 trajectory  $k$ , we define  $R(u, k)$  as follows:

- 262 •  $R(u, k) = \text{SUCCESS}$  if  $k[-1][1] = \text{goal}$
- 263 •  $R(u, k) = \text{DEADEND}$  if  $k[-1][1] \in P$
- 264 •  $R(u, k) = \text{MAX STEP}$  if  $|k| > \text{max}$
- 265 •  $R(u, k) = \text{INVALID}$  if  $\exists a \in k[:][0] \mid a \notin A$

266 <sup>2</sup> In the paper, we set the maximum steps  $max = 200$  for each environment. We denote  $P$  as the set  
267 of deadend states, i.e. set of pits and walls in the grid, and all the states outside the grid.

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

## 272 D.2 Cognitive map design choices

273 Our construction of the cognitive map is straightforward; We sample all 4 possible directions("up",  
274 "down", "right", and "left") as actions, and we propagate each action iteratively until reaching the goal  
275 state. After reaching the goal, we backtrack until reaching the start state. For ablation experiments, we  
276 denote ALL as sampling all 4 actions and BACKTRACK as appending backtrack trace after reaching  
277 the goal. We experiment effects of excluding each component when designing the cognitive map.  
278 First, we see the effect of excluding the backtrack stage(denoted as "w.o. BACKTRACK"). Also, we  
279 observe the effect when we only sample possible moves instead of all moves(denoted as "w.o. ALL").  
280 See Appendix E.1 , E.2 for examples.

281 **Marking deadend** In this paper, we variate two different cognitive maps when marking the deadend  
282 state. We can either verbalize all samples and mark the actions resulting deadend state with a special  
283 token(<deadend>)(denoted as "MARKING deadend"), or just verbalize all the actions including the  
284 deadend states(denoted as "w.o. MARKING"). For example, the verbalized cognitive map at state (0,  
285 0) in Figure 2 is as follows:

286 MARKING deadend: "(0, 1) up (0, -1) <deadend> (1, 0) right (-1, 0) <deadend>"  
287 w.o. MARKING: (0, 1) up (0, -1) down (1, 0) right (-1, 0) left".  
288 See Appendix E.1 for the whole construction.

289 **Backward cognitive map construction: BWD** Backward chaining is a powerful approach that  
290 simplifies complex problems by focusing on the desired outcome and systematically working back to  
291 the starting point. LAMBADA (Kazemi et al., 2023) shows that backward chaining helps in reasoning  
292 tasks. We adopt the intuition to see if constructing the cognitive map in a backward manner enhances  
293 the planning ability of the language model.

294 In this paper, we define two types of construction, FWD and BWD. For FWD, the construction from the  
295 start to the goal is identical to the procedure stated in Appendix C. For BWD, we build the reversed  
296 cognitive map starting from the goal state to the start state. We have the sampling function identical  
297 to FWD( $S$ ). The main experiments are set to BWD. For the propagation stage, we expect a reverse  
298 transition state  $T^{-1}(s, a) \in \mathcal{S}$  for a given state  $s$  and action  $a$ . For the backtracking stage, we  
299 backtrack the path from start to goal by iteratively searching  $T(s, a) \in \mathcal{S}$ . See Appendix E.2 for the  
300 example.

---

<sup>2</sup>Although our main analysis consists only of success rate, we provide plots for different types of fails in Appendix F.3

301 **E Cognitive map description**

302 **E.1 Cognitive map example: FWD**

303 Table 3 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 3: FWD cognitive map example for each experiment

304 **E.2 Cognitive map example: BWD**

305 Table 4 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)\nrigh\tn(1, 0)\nrigh\tn(2, 0)\nup\tn(2, 1)\nrigh\tn(3, 1)\nup\tn(3, 2)
w.o. ALL BACKTRACK	Thought:\nStep 1:\n(3, 1)\nup\tnStep 2:\n(2, 1)\nrigh\tnStep 3:\n(2, 0)\nup\tnStep 4:\n(1, 0)\nrigh\tnStep 5:\n(0, 0)\nrigh
w.o. ALL	Thought:\nStep 1:\n(3, 1)\nup\tnStep 2:\n(2, 1)\nrigh\tnStep 3:\n(2, 0)\nup\tnStep 4:\n(1, 0)\nrigh\tnStep 5:\n(0, 0)\nrigh\tnBacktrack:\n(0, 0)\nrigh\tn(1, 0)\nrigh\tn(2, 0)\nup\tn(2, 1)\nrigh\tn(3, 1)\nup\tn(3, 2)
w.o. MARKING BACKTRACK	Thought:\nStep 1:\n(3, 3)\ndown\tn(3, 1)\nup\tn(2, 2)\nrigh\tn(4, 2)\nleft\tnStep 2:\n(3, 2)\ndown\tn(3, 0)\nup\tn(2, 1)\nrigh\tn(4, 1)\nleft\tnStep 3:\n(2, 2)\ndown\tn(2, 0)\nup\tn(1, 1)\nrigh\tn(3, 1)\nleft\tnStep 4:\n(2, 1)\ndown\tn(2, -1)\nup\tn(1, 0)\nrigh\tn(3, 0)\nleft\tnStep 5:\n(1, 1)\ndown\tn(1, -1)\nup\tn(0, 0)\nrigh\tn(2, 0)\nleft
w.o. MARKING	Thought:\nStep 1:\n(3, 3)\ndown\tn(3, 1)\nup\tn(2, 2)\nrigh\tn(4, 2)\nleft\tnStep 2:\n(3, 2)\ndown\tn(3, 0)\nup\tn(2, 1)\nrigh\tn(4, 1)\nleft\tnStep 3:\n(2, 2)\ndown\tn(2, 0)\nup\tn(1, 1)\nrigh\tn(3, 1)\nleft\tnStep 4:\n(2, 1)\ndown\tn(2, -1)\nup\tn(1, 0)\nrigh\tn(3, 0)\nleft\tnStep 5:\n(1, 1)\ndown\tn(1, -1)\nup\tn(0, 0)\nrigh\tn(2, 0)\nleft\tnBacktrack:\n(0, 0)\nrigh\tn(1, 0)\nrigh\tn(2, 0)\nup\tn(2, 1)\nrigh\tn(3, 1)\nup\tn(3, 2)
w.o. BACKTRACK	Thought:\nStep 1:\n(3, 3)\ncut\tn(3, 1)\nup\tn(2, 2)\ncut\tn(4, 2)\ncut\tnStep 2:\n(3, 2)\ncut\tn(3, 0)\ncut\tn(2, 1)\nrigh\tn(4, 1)\ncut\tnStep 3:\n(2, 2)\ncut\tn(2, 0)\nup\tn(1, 1)\ncut\tn(3, 1)\ncut\tnStep 4:\n(2, 1)\ncut\tn(2, -1)\ncut\tn(1, 0)\nrigh\tn(3, 0)\ncut\tnStep 5:\n(1, 1)\ncut\tn(1, -1)\ncut\tn(0, 0)\nrigh\tn(2, 0)\ncut
MARKING deadend	Thought:\nStep 1:\n(3, 3)\ncut\tn(3, 1)\nup\tn(2, 2)\ncut\tn(4, 2)\ncut\tnStep 2:\n(3, 2)\ncut\tn(3, 0)\ncut\tn(2, 1)\nrigh\tn(4, 1)\ncut\tnStep 3:\n(2, 2)\ncut\tn(2, 0)\nup\tn(1, 1)\ncut\tn(3, 1)\ncut\tnStep 4:\n(2, 1)\ncut\tn(2, -1)\ncut\tn(1, 0)\nrigh\tn(3, 0)\ncut\tnStep 5:\n(1, 1)\ncut\tn(1, -1)\ncut\tn(0, 0)\nrigh\tn(2, 0)\ncut\tnBacktrack:\n(0, 0)\nrigh\tn(1, 0)\nrigh\tn(2, 0)\nup\tn(2, 1)\nrigh\tn(3, 1)\nup\tn(3, 2)

Table 4: BWD cognitive map example for each experiment



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

Table 5: 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 E for actual prompts.

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

Table 6: 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	<b>0.765</b>	0.705
FWD	0.528	0.516	<b>0.618</b>	0.585
Reachable	w.o. MARKING BACKTRACK	w.o. BACKTRACK	w.o. MARKING	MARKING deadend
BWD	0.739	0.852	0.724	<b>0.885</b>
FWD	0.672	0.624	0.816	<b>0.854</b>

Table 7: Planning performance with BACKTRACK exclusion

## 306 F Additional results

### 307 F.1 Additional investigations

308 **Cognitive map compared with baselines** Table 5 shows both the optimal and reachable rate  
309 of plans generated by each experiment. As shown in the table, w.o. MARKING shows the best  
310 performance among the experiments both for BWD and FWD cognitive map construction in the  
311 optimal planning. (76.5% and 61.8% for BWD and FWD construction, respectively). Unlike in the  
312 optimal planning setting, MARKING deadend shows the best performance both for BWD and FWD  
313 cognitive map construction(88.5% and 85.4% for BWD and FWD construction, respectively). **Our**  
314 **experiments show that cognitive maps improve performance in the Gridworld path planning**  
315 **task.** It boosts the optimal planning performance by up to 57.5% and reachable planning by up to  
316 56.4% with implicit fine-tuning baseline, and enhancing optimal planning by up to 50% and reachable  
317 planning by up to 54.6% with CoT fine-tuning baseline. We also qualitatively show the planning  
318 performance in both optimal and reachable settings with respect to the size of the grid in Table 1.

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

325 **BWD approach enhances the performance of the cognitive map** As shown in Table 5, both the  
326 optimal and reachable planning got the highest performance with BWD cognitive map construction.  
327 This aligns with the findings in LAMBADA (Kazemi et al., 2023), where they show that backward  
328 chaining helps in reasoning tasks.

329 **Effect of ALL inclusion** Both cognitive map constructions w.o. ALL (w.o. ALL BACKTRACK:  
330 29.6% and 29.5% for BWD and FWD construction, respectively; w.o. ALL: 27.7% and 29.0% for  
331 BWD and FWD construction, respectively) suffer at generating the optimal plan, achieving under 30%  
332 success rate. However, every cognitive map construction with the inclusion of ALL shows better  
333 performance, with FWD MARKING deadend being the lowest among them(58.5%).

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

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

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

### 348 F.2 Vizualization for optimal planning experiments

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

351 **FWD construction** See Figure 4 for success rate of each experiment.

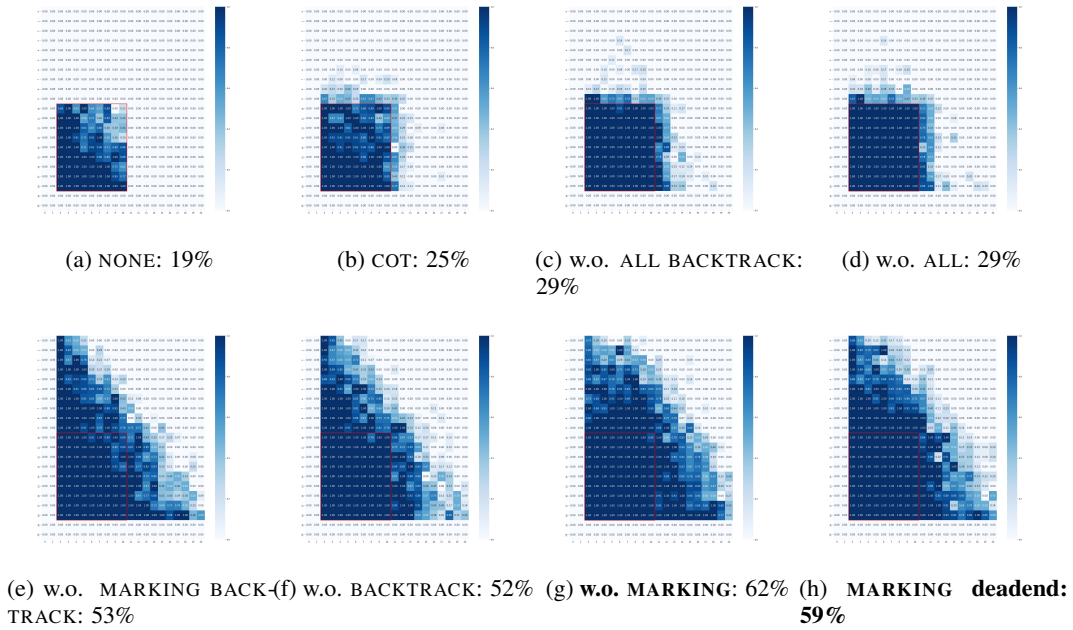


Figure 4: Success rate for optimal planning, FWD construction

352 **BWD construction** See Figure 5 for success rate of each experiment.

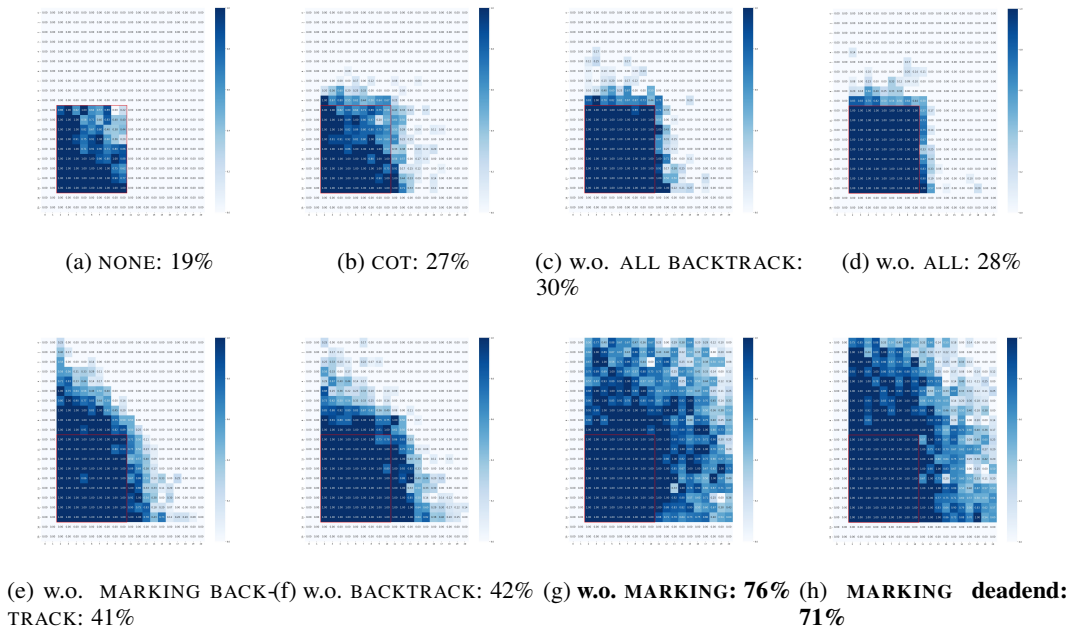


Figure 5: Success rate for optimal planning, BWD construction

### 353 F.3 Visualization for reachable planning experiments

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

356 **FWD construction** For success rate, see Figure 6. For failure cases, see Figure 7 for deadend,  
 357 Figure 8 for max step, and Figure 9 for invalid rate.

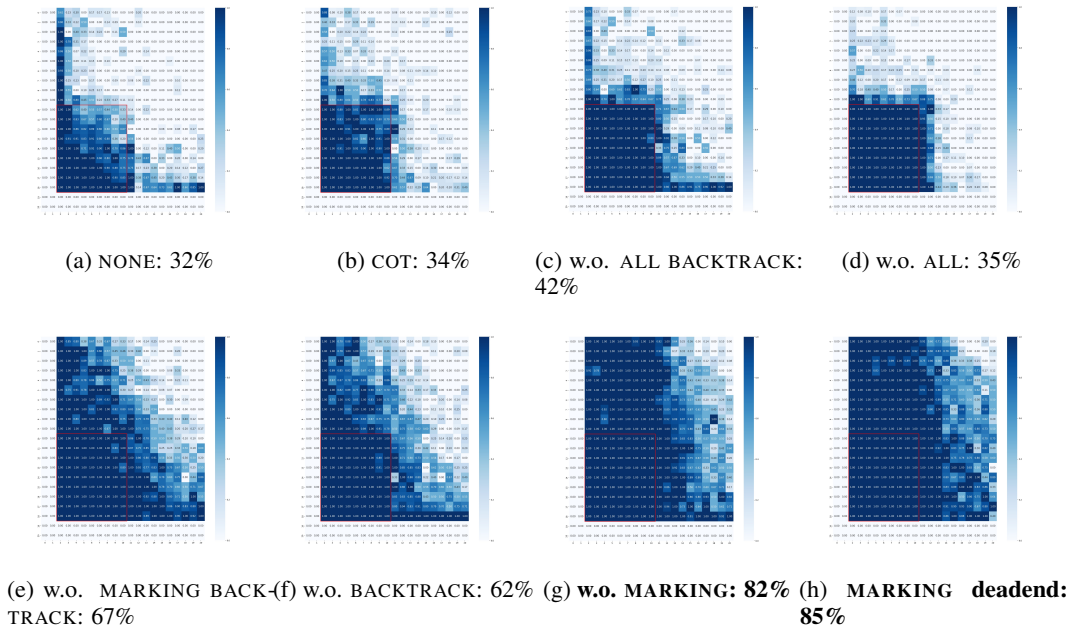


Figure 6: Success rate for reachable planning, FWD construction

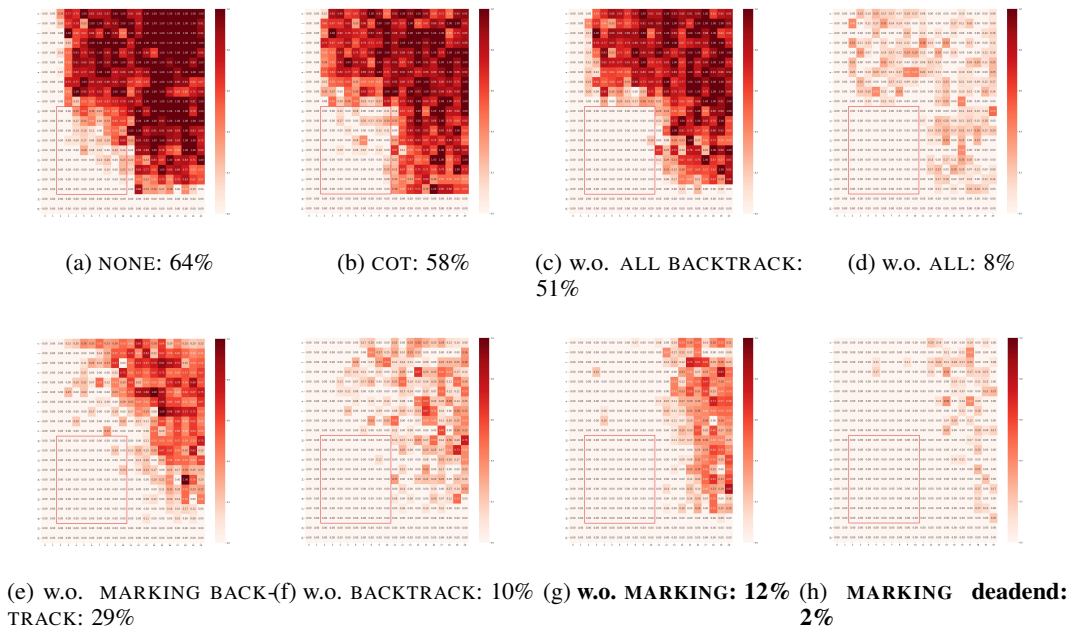


Figure 7: Deadend rate for reachable planning, FWD construction

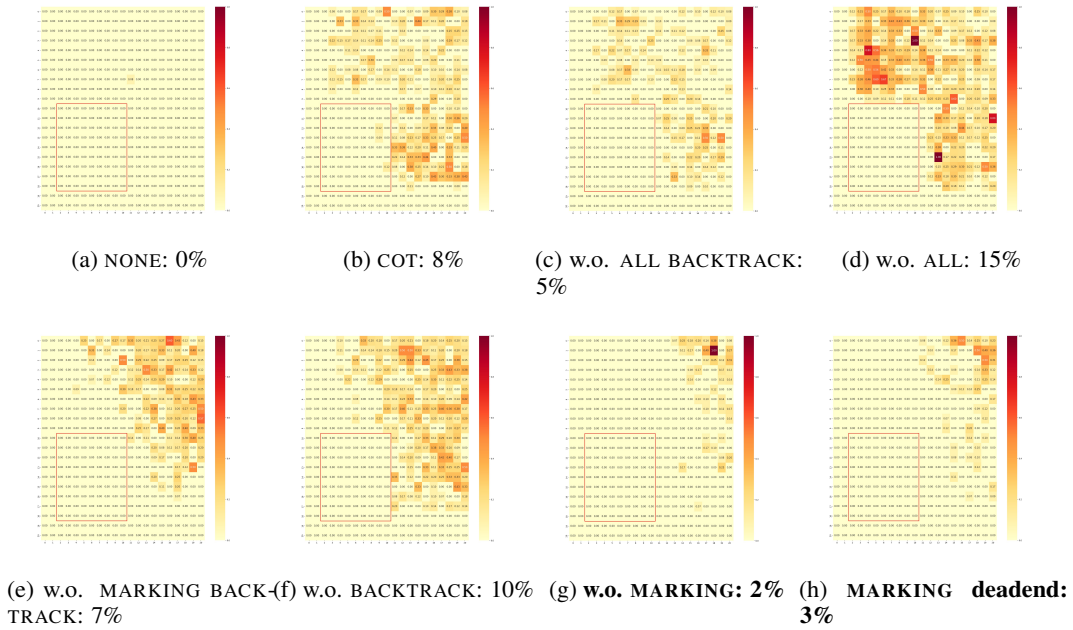


Figure 8: Max step rate for reachable planning, FWD construction

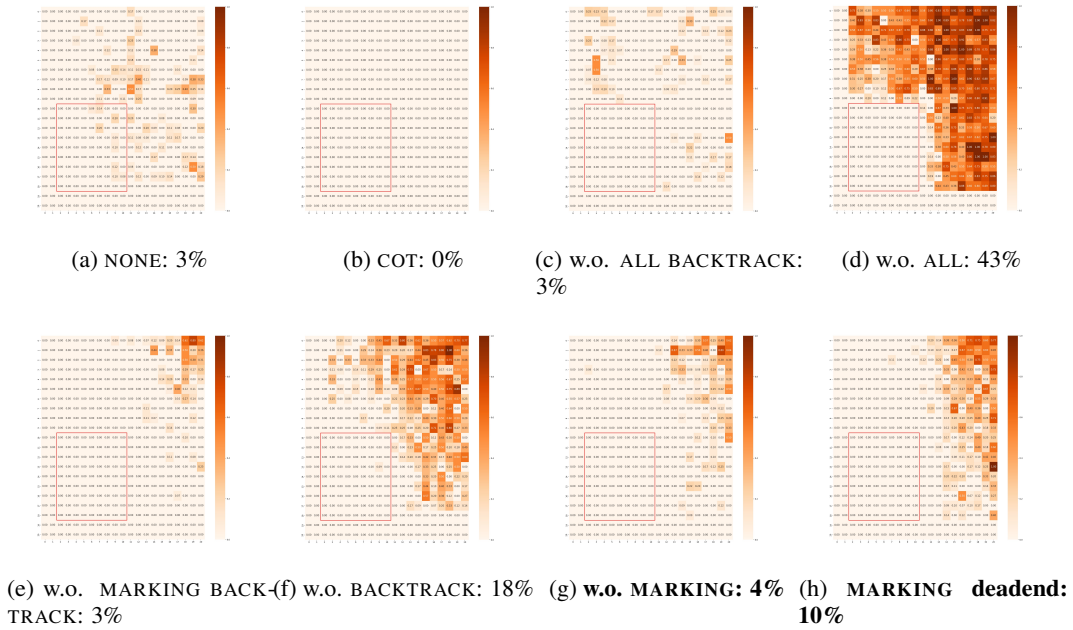


Figure 9: Invalid rate for reachable planning, FWD construction

358 **BWD construction** For success rate, see Figure 10. For failure cases, see Figure 11 for deadend,  
 359 Figure 12 for max step, and Figure 13 for invalid rate.

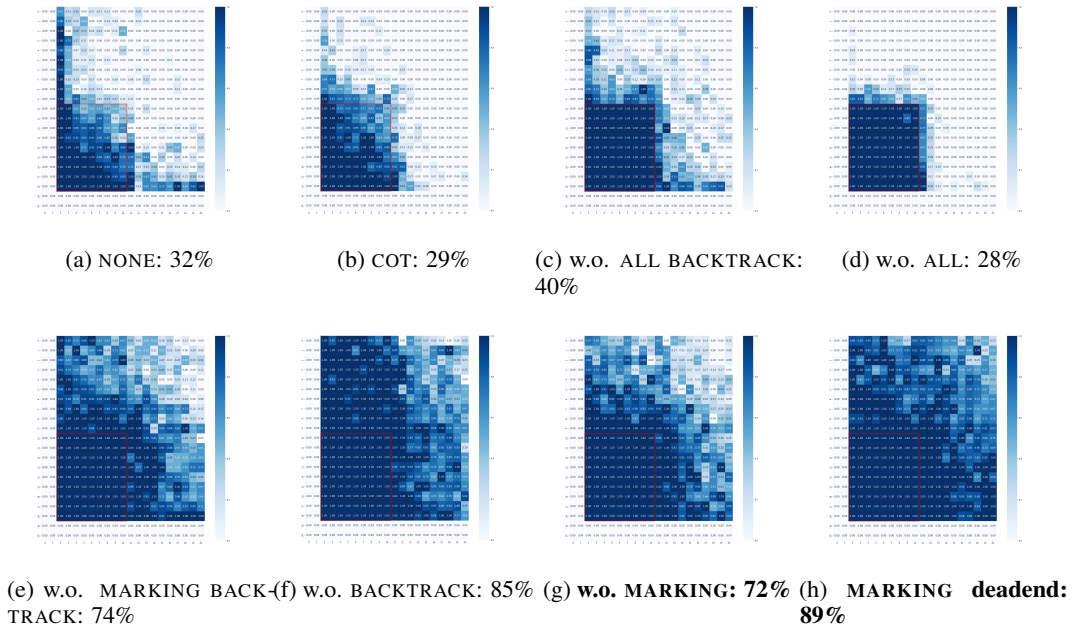


Figure 10: Success rate for reachable planning, BWD construction

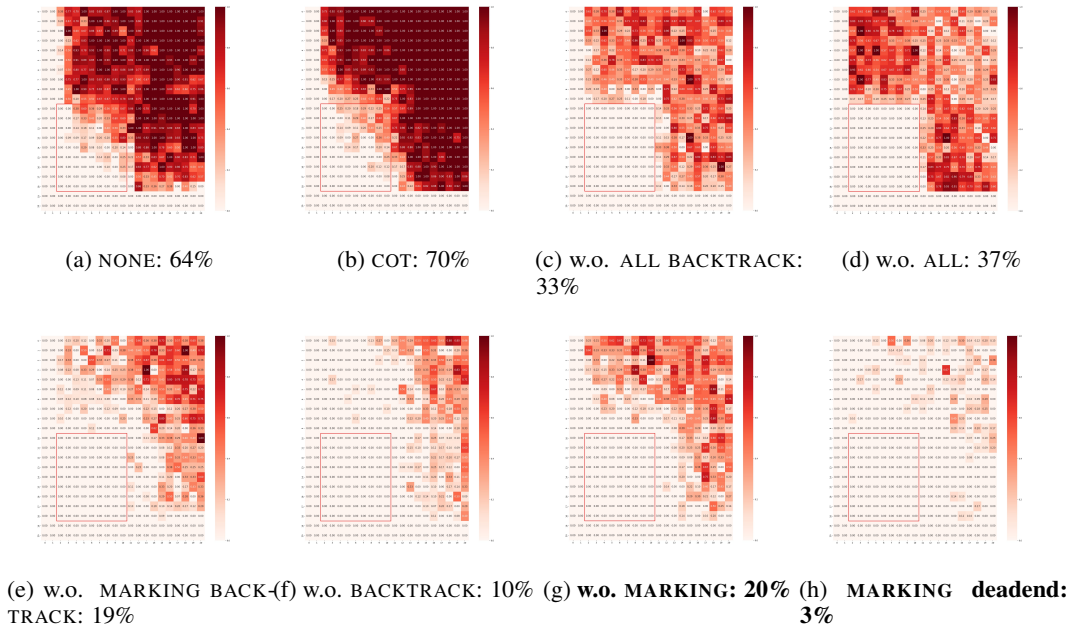


Figure 11: Deadend rate for reachable planning, BWD construction

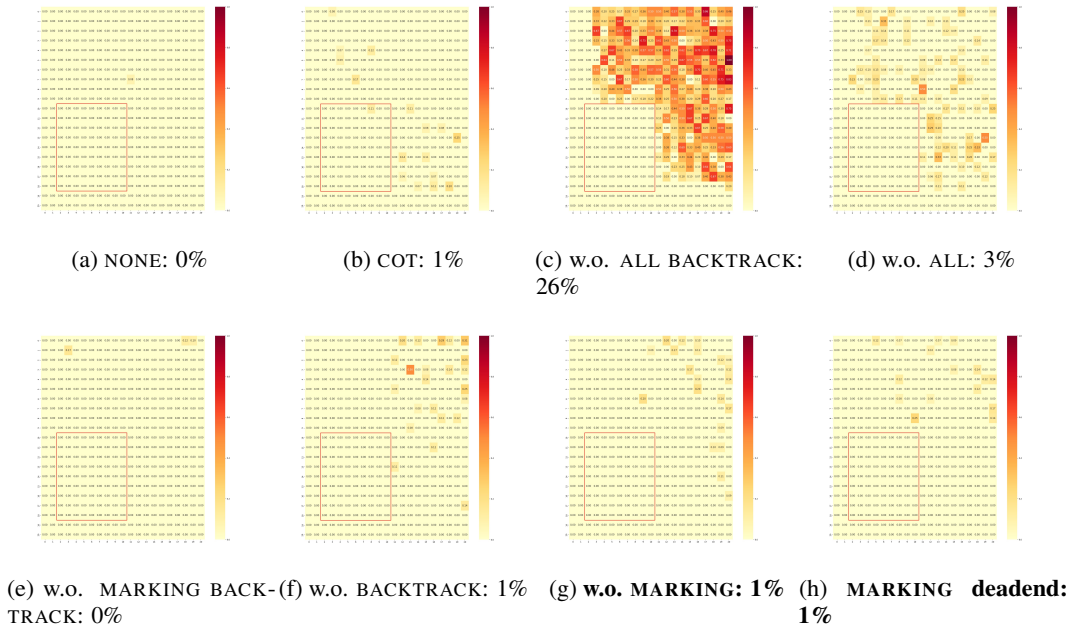


Figure 12: Max step rate for reachable planning, BWD construction

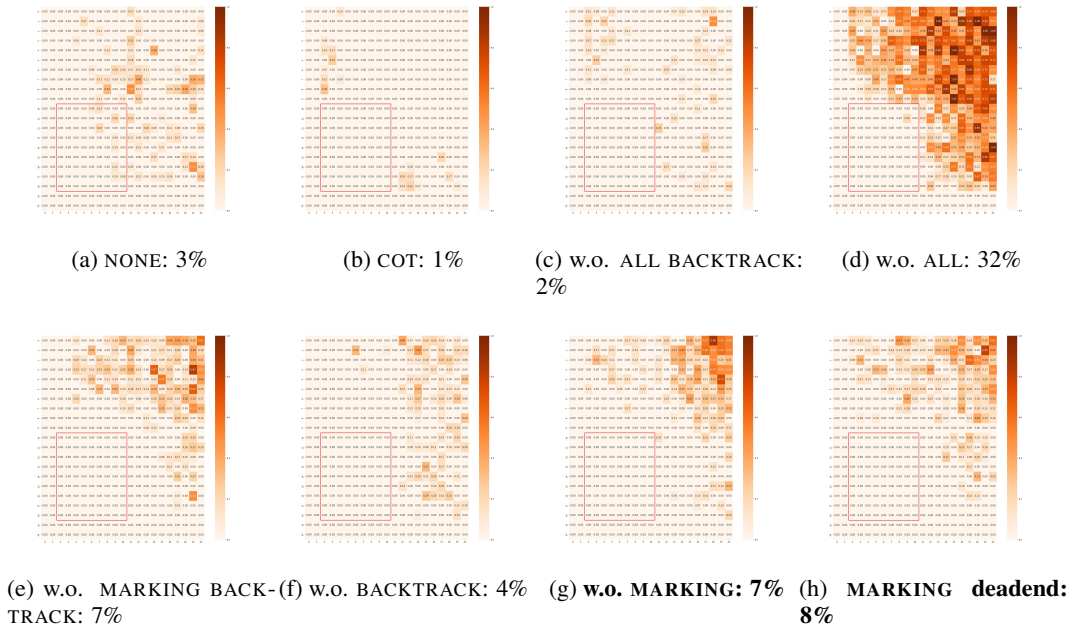


Figure 13: Invalid rate for reachable planning, BWD construction