

# DeepMaze: A Maze Benchmark for Quantifying Reasoning Depth in Language Models

Anonymous ACL submission

## Abstract

Existing benchmarks for evaluating reasoning in large language models primarily emphasize final-answer correctness, making it difficult to distinguish genuine multi-step reasoning from statistical shortcuts and prior knowledge exploitation. Accurately measuring reasoning depth requires environments where this confounding is eliminated; however, existing benchmarks retain data contamination and static test distribution bias that enable prior knowledge to masquerade as reasoning, preventing clean capability isolation. To overcome this, we introduce DEEPMAZE, a minimalist benchmark of procedurally generated environments with rigorously controlled topology. Its dual-task architecture—comprising *planning* under full observability and *exploration* under partial observability—inherently necessitates sustained, state-consistent reasoning by requiring models to dynamically track environmental states across sequential actions. Within this environment, we define a reasoning depth metric that quantifies the length of state-consistent action sequences, explicitly decoupling process quality from outcome success. This design isolates LLMs’ core reasoning capabilities under controlled conditions, establishing a foundation for evaluating their true multi-step reasoning proficiency independent of domain-specific knowledge or outcome-driven shortcuts.

## 1 Introduction

Spatial reasoning ability refers to the capacity to construct, organize, manipulate, and transform spatial information and spatial relations (Newcombe and Shipley, 2014). In real-world and embodied tasks, this capability requires sustained, state-consistent reasoning over multi-step actions to interact effectively with dynamic environments. In recent years, a growing number of benchmarks have been proposed to evaluate the spatial reasoning capabilities of large language models and

multimodal models, providing important tools for understanding models’ spatial cognition. However, existing approaches still suffer from several critical limitations, which constrain their ability to faithfully characterize models’ genuine reasoning capabilities.

First, existing benchmarks fail to disentangle and explicitly characterize reasoning depth. Reasoning depth requires quantifying the length of state-consistent action sequences—where each decision depends on dynamically updated internal states formed from prior actions—in ways that cannot be replaced by statistical shortcuts. Crucially, most evaluations equate reasoning with final-answer correctness, conflating genuine multi-step reasoning with outcome success. This prevents independent assessment of how much sustained, state-consistent reasoning a model performs during execution, as even memorized solutions yield correct answers without intermediate state updates.

Second, evaluation outcomes are heavily confounded by static test distribution bias. Fixed or limited-scale datasets contain inherent statistical regularities (e.g., frequent path patterns or layout symmetries), enabling models to exploit pattern matching and prior knowledge instead of explicit reasoning. Models succeed via statistical shortcuts that bypass multi-step state-dependent decisions, as environmental complexity alone cannot eliminate these biases when test distributions remain static. Consequently, performance reflects dataset memorization rather than adaptive reasoning capability.

Third, data contamination critically undermines evaluation validity. When test instances leak into pretraining or fine-tuning data—a common issue in benchmarks with limited procedural generation—models leverage prior knowledge to masquerade as reasoning. Performance then measures memorization of specific instances rather than state-consistent reasoning in novel situations, severely

085 weakening the benchmark’s ability to isolate  
086 true reasoning competence. This contamination,  
087 combined with static distributions, creates an  
088 environment where spurious shortcuts dominate  
089 genuine capability assessment.

090 Guided by the principle of “less is more”—  
091 namely, isolating the core mechanisms of reason-  
092 ing and memory by stripping away semantic con-  
093 tent and prior-knowledge cues through a minimal-  
094 ist maze environment—we introduce DEEPMAZE,  
095 a benchmark built upon procedurally generated  
096 mazes. DEEPMAZE adopts a two-level progressive  
097 task hierarchy to systematically decompose and  
098 probe different dimensions of spatial reasoning  
099 ability: (1) **Planning**: given the full observation of  
100 maze layout, the model plans and executes moves  
101 step-by-step—with multiple attempts allowed—  
102 until reaching the goal, evaluating its long-horizon  
103 planning capability and enabling the assessment of  
104 reasoning depth; (2) **Exploration**: the model is  
105 restricted to local observations of its surrounding  
106 environment and must make decisions while  
107 exploring, thereby forcing it to construct and  
108 continuously update an internal spatial representa-  
109 tion; Based on the above design philosophy and  
110 methodological considerations, this paper makes  
111 the following main contributions:

112 (1) We propose DEEPMAZE, a minimalist bench-  
113 mark leveraging procedurally generated environ-  
114 ments that isolates core reasoning capabilities from  
115 prior knowledge exploitation. The benchmark  
116 features a rigorously controlled dataset of over  
117 50,000 distinct maze instances spanning 10 grid  
118 scales (from  $5 \times 5$  to  $29 \times 29$ ), where topological  
119 diversity eliminates statistical shortcuts while  
120 enhancing generalization.

121 (2) We introduce a *Reasoning Depth Metric*  
122 that quantifies sustained state-consistent reasoning  
123 through consecutive valid actions, decoupling  
124 process quality from outcome success. This metric  
125 provides the first standardized method to measure  
126 reasoning chain length independent of final-answer  
127 correctness.

128 (3) We establish a dual-task evaluation framework  
129 —comprising *planning* (full observability) and  
130 *exploration* (partial observability)—that system-  
131 atically probes state-tracking capabilities under  
132 increasing environmental complexity. Compre-  
133 hensive analyses across state-of-the-art models  
134 reveal a fundamental *reasoning gap*: high success  
135 rates often mask failures in maintaining state-  
136 consistent logic as complexity grows.

## 2 Related Works 137

### 2.1 General Reasoning Benchmarks 138

139 Evaluating multi-step reasoning in LLMs has led  
140 to diverse benchmarks spanning mathematical,  
141 logical, and physical domains. Mathematical  
142 benchmarks like MATH (Hendrycks et al., 2021)  
143 and GSM8K (Cobbe et al., 2021) assess solution  
144 correctness but **fail to disentangle reasoning  
145 depth from formulaic shortcuts**. For example, a  
146 model may solve a geometry problem by directly  
147 recalling “area =  $\pi r^2$ ” without reconstructing  
148 the derivation—compressing multi-step spatial  
149 reasoning into a single memorized step. Sim-  
150 ilarly, theorem-proving benchmarks (e.g., Lean  
151 (Yang et al., 2023)) rely on formal syntax that  
152 allows proof-skipping via lemma reuse. Logical  
153 reasoning datasets (e.g., ReClor (Yu et al., 2020),  
154 LOGIQA (Liu et al., 2020)) emphasize deductive  
155 validity but lack metrics for intermediate state con-  
156 sistency, conflating correct conclusions with sound  
157 inference processes. Crucially, these benchmarks  
158 share a core flaw: static test distributions enable  
159 exploitation of secondary conclusions, preventing  
160 isolation of pure reasoning depth. Our work  
161 bridges this gap by designing a domain-agnostic  
162 evaluation framework where reasoning depth is  
163 inherently uncompressible due to procedural topol-  
164 ogy control.

### 2.2 Spatial Reasoning Benchmarks 165

166 Existing spatial benchmarks (e.g., SpatialBench  
167 (Xu et al., 2025), ALFWorld (Shridhar et al.,  
168 2021)) inherit the above limitations. Real-  
169 world scene datasets embed semantic priors (e.g.,  
170 “kitchens contain fridges”), allowing models to by-  
171 pass geometric reasoning via object co-occurrence  
172 statistics. Even procedurally generated environ-  
173 ments (e.g., VirtualHome (Puig et al., 2018)) suffer  
174 from layout regularities that enable memorization.  
175 Critically, **none decouple reasoning depth from  
176 outcome success**—a gap our state-consistent met-  
177 ric directly addresses.

### 2.3 Maze-based Benchmarks and Studies 178

179 Maze tasks have long served as a classical  
180 paradigm for evaluating multi-step reasoning due  
181 to their simple structure, abstract semantics, and in-  
182 herent requirement for sequential decision-making.  
183 In reinforcement learning literature, benchmarks  
184 such as MiniGrid (Chevalier-Boisvert et al., 2023)  
185 and BabyAI (Chevalier-Boisvert et al., 2019)

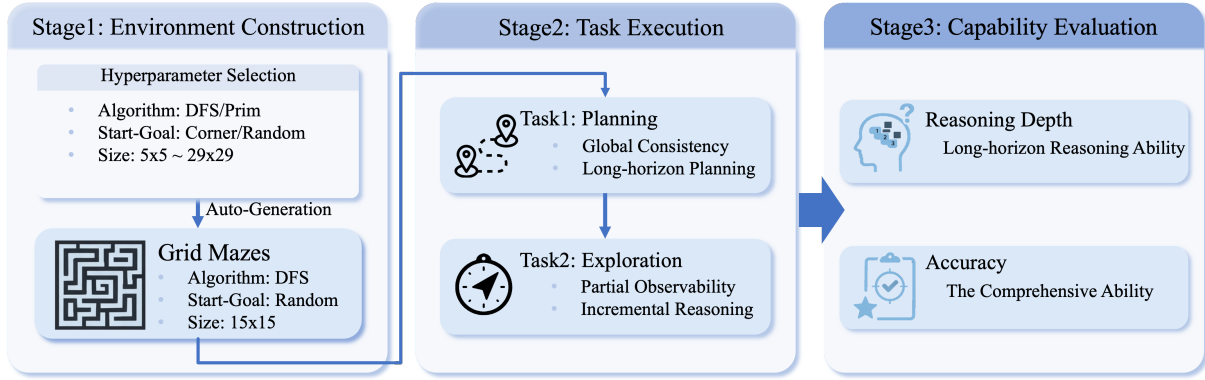


Figure 1: Overview of DeepMaze

186 have been widely adopted to assess navigation  
 187 capabilities. However, these frameworks primarily  
 188 target RL agents and lack standardized metrics for  
 189 evaluating reasoning depth in language models.

190 Einarsson (Einarsson, 2025) introduced  
 191 MazeEval to evaluate LLM decision-making in  
 192 coordinate-based mazes using sparse feedback  
 193 signals, with findings showing significant  
 194 performance degradation in larger mazes and  
 195 non-English contexts. Dao et al. (Dao and Vu,  
 196 2025) developed AlphaMaze with tokenized  
 197 maze representations trained via SFT and  
 198 GRPO, but their approach suffers from quadratic  
 199 sequence growth with maze size and lacks partial  
 200 observability support, limiting scalability for  
 201 complex reasoning evaluation.

202 Critically, existing maze-based studies predomi-  
 203 nantly treat mazes as *end tasks* rather than *reason-*  
 204 *ing probes*. They focus on navigation success rates  
 205 within the maze domain without establishing how  
 206 maze performance reflects fundamental reasoning  
 207 capabilities such as state consistency maintenance  
 208 or resistance to statistical shortcuts. These works  
 209 fail to isolate and quantify the core reasoning  
 210 mechanisms that transcend specific maze configu-  
 211 rations.

212 In contrast, our benchmark treats mazes as a  
 213 *minimalist reasoning microscope*. By stripping  
 214 away semantic content and prior-knowledge cues  
 215 through controlled procedural generation, Deep-  
 216 Maze directly measures state-consistent reasoning  
 217 depth—quantified by consecutive valid actions—  
 218 across two complementary conditions: planning  
 219 under full observability and exploration under  
 220 partial observability. This design enables precise  
 221 characterization of how reasoning quality degrades  
 222 with environmental complexity, providing transfer-  
 223 able insights about model reasoning capabilities

beyond maze navigation.

### 3 Benchmark Description

224 DEEPMAZE defines three progressively challeng-  
 225 ing tasks to isolate reasoning capabilities. Built  
 226 on procedural maze generation, it enables rigorous  
 227 control over topology and observation constraints.  
 228  
 229

#### 3.1 Environment Formalization

230 A maze is a grid graph  $\mathcal{M} = (V, E)$  where  $V \subseteq$   
 231  $[H] \times [W]$  contains traversable cells. State  $s_t =$   
 232  $(p_t, g)$  tracks agent position  $p_t$  and goal  $g$ . Actions  
 233  $\mathcal{A} = \{U, D, L, R\}$  trigger deterministic transitions:  
 234

$$s_{t+1} = \mathcal{T}(s_t, a_t) = \begin{cases} (p_t + \Delta a_t, g) & \text{if valid} \\ s_t & \text{otherwise} \end{cases}$$

235 Observations depend on radius  $r$ :  $o_t =$   
 236  $\mathcal{O}(p_t, g; \mathcal{M}, r)$ .  
 237

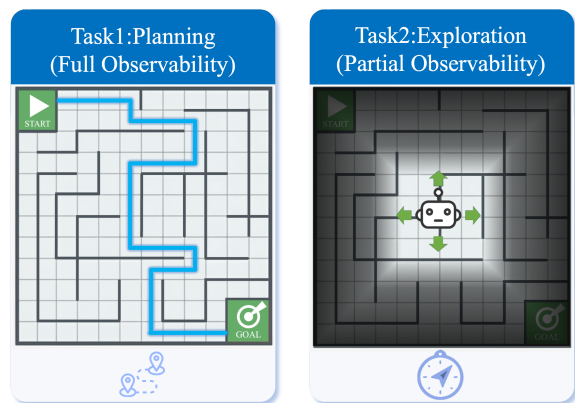


Figure 2: DeepMaze’s three-task framework: (Left) Global planning with full observability; (Right) Exploration under partial observability.

#### 3.2 Task 1: Global Planning

238 The agent receives full maze observation  $o_t^{\text{global}}$   
 239 and outputs a single action sequence  $\mathbf{a}_{t:t+K}$ .  
 240

Table 1: Comparison of Maze-Based Benchmarks for LLM Reasoning Evaluation

Evaluation Dimension	MazeEval (Einarsson, 2025)	MazeBench (Dao and Vu, 2025)	DeepMaze (Ours)
Reasoning depth quantification	✗	✗	✓
State-consistency measurement	✗	✗	✓
Observation modes supported	Partial only	Global only	<b>Global &amp; Partial</b>
Maze size range	5×5–15×15	5×5	<b>5×5–29×29</b>
Evaluation dimensions	Outcome-only	Outcome-only	<b>Process &amp; Outcome</b>
Token efficiency (global)	N/A	≈ 32n <sup>2</sup>	<b>n<sup>2</sup></b>
Primary purpose	Navigation task	Navigation task	<b>Reasoning probe</b>
Task paradigm	Exploration	Planning	<b>Planning &amp; Exploration</b>

Execution follows strict all-or-nothing rule:

$$\mathcal{E}(s_t, \mathbf{a}) = \begin{cases} s_t & \text{if any invalid action} \\ s_{t+K} & \text{otherwise} \end{cases}$$

Unlike prior work, we evaluate plan coherence beyond endpoint success.

### 3.3 Task 2: Local Exploration

The agent observes only a local window of radius  $r \ll \min(H, W)$ :

$$o_t = \mathcal{O}(p_t, g; \mathcal{M}, r).$$

It outputs short action sequences incrementally based on history  $o_{\leq t}$ , evaluating online reasoning under uncertainty.

### 3.4 Evaluation Metrics

For  $N$  episodes  $\mathcal{E}$  with max steps  $T_{\max}$ :

#### Success Rate

$$\mathbb{I}_{\text{succ}}(e) = \begin{cases} 1 & \text{if reached } g \text{ within } T_{\max} \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Success} = \frac{1}{N} \sum_{e \in \mathcal{E}} \mathbb{I}_{\text{succ}}(e).$$

**Reasoning Depth** For episode  $e$  with  $m_e$  planning attempts:

$$k_{e,i} = \max\{k \mid \forall j \leq k: a_j^{(i)} \text{ valid}\}$$

$$\text{RD}(e) = \frac{1}{m_e} \sum_{i=1}^{m_e} k_{e,i}$$

$$\text{Reasoning Depth} = \frac{1}{N} \sum_{e \in \mathcal{E}} \text{RD}(e)$$

## 4 Experiments

### 4.1 Maze Configuration

We generate mazes using the *Recursive Division Algorithm*, producing connected acyclic structures in ASCII format (walls ‘#’, empty cells ‘0’, agent ‘A’, goal ‘G’). Maze sizes are scaled continuously:

- Task 1: 5 × 5 to 29 × 29 grids (critical thresholds at 11 × 11, 13 × 13)
- Task 2: 5 × 5 to 9 × 9 grids with partial observability

For Task 2, we implement configurable visibility radii (1-7 cells) using ray-casting for occlusion simulation, with 3-cell radius as default. All experiments use 1,000 trials per condition (600 for Qwen-3-max-preview due to computational constraints).

### 4.2 Models and Context Mechanism

We evaluate the following models with DeepSeek-V3.2 as our primary baseline:

- **DeepSeek-V3.2** (baseline): Standard variant with temperature=0 and 1.0 settings
- **DeepSeek-V3.2-Thinking**: Enhanced reasoning variant with iterative problem-solving capability
- **Proprietary models**: GPT-4o, GPT-5.1, GPT-o3
- **Other open-source**: Qwen-3-max-preview

289 All models operate with a context window 339  
290 containing the **most recent  $k$  interaction cycles**, 340  
291 where each cycle consists of an *action sequence* 341  
292 proposed by the model followed by *execution* 342  
293 *feedback* from the environment (valid moves, 343  
294 collisions, and partial observations). This forms 344  
295 a closed-loop interaction history that enables state 345  
296 tracking and adaptive decision-making. For Task 346  
297 2, we systematically evaluate different context 347  
298 window sizes ( $k \in \{3, 5, 7\}$ ) to quantify memory 348  
299 requirements for effective exploration. 349

300 **Task 2 constraints:** Step limits ( $k \times n^2$ ,  $k \in$  350  
301  $\{1.0, 2.0\}$ ), exploration patience ( $m \in \{3, 5, 7, 9\}$  351  
302 consecutive non-exploring steps), memory win- 352  
303 dow sizes (3-7 interaction cycles), and visibility 353  
304 radii (1-7 cells). Section 3.3 presents our ablation 354  
305 study on memory window size, revealing critical 355  
306 thresholds for maintaining state consistency during 356  
307 exploration. 357

## 308 4.3 Results and Analysis 358

### 309 4.3.1 Task 1 359

310 **Accuracy trends.** The experimental results re- 360  
311 veal critical insights into model performance 361  
312 under increasing spatial complexity. As shown 362  
313 in Figure 3, all non-reasoning models exhibit 363  
314 a clear negative correlation between maze size 364  
315 and success rate, with degradation becoming 365  
316 pronounced beyond  $11 \times 11$  grids. However, the 366  
317 rate and severity of decline vary dramatically 367  
318 across architectures and configurations. 368

319 The *DeepSeek-V3.2-Thinking* variant demon- 370  
320 strates exceptional robustness, maintaining perfect 371  
321 success rates (1.000) across all grid sizes from 372  
322  $5 \times 5$  to  $29 \times 29$ . This suggests that its “thinking” 373  
323 mechanism—enabling multi-step reasoning and 374  
324 iterative planning—effectively mitigates combina- 375  
325 torial complexity in larger mazes. In contrast, the 376  
326 non-reasoning *DeepSeek-V3.2 (temperature=0)* 377  
327 configuration suffers a sharp performance cliff: 378  
328 while achieving 0.80 success at  $5 \times 5$ , it drops 379  
329 to 0.00 for all grids  $\geq 13 \times 13$ . This abrupt 380  
330 failure indicates that greedy decoding without 381  
331 stochastic exploration cannot handle pathfinding 382  
332 tasks requiring long-horizon planning. 383

333 The divergence in medium-to-large grids ( $\geq$  384  
334  $15 \times 15$ ) is particularly instructive: while all non- 385  
335 reasoning models fail completely in this regime, 386  
336 the “thinking” variant sustains perfect perfor- 387  
337 mance. This underscores the critical role of 388  
338 explicit reasoning mechanisms in solving scalable 389

spatial reasoning tasks. The results align with 339  
theoretical expectations that maze-solving requires 340  
hierarchical planning—a capability embedded in 341  
the “thinking” configuration but absent in standard 342  
autoregressive decoding. 343

344 These findings have practical implications for 345  
346 real-world navigation tasks where environmental 347  
348 complexity is unbounded. The stark contrast 349  
350 between reasoning and non-reasoning configu- 351  
352 rations suggests future work should prioritize 353  
354 architectures with built-in planning capabilities 355

356 **Reasoning depth analysis.** The grouped bar 357  
358 chart in Figure 4 reveals a fundamental dichotomy 359  
360 in problem-solving approaches between reasoning- 361  
362 enabled and standard architectures. At the  $7 \times 7$  363  
364 scale, all models exhibit comparable reasoning 365  
366 depth, confirming that small-scale mazes do not 367  
368 adequately stress-test planning capabilities. How- 369  
370 ever, as complexity increases to  $15 \times 15$  and  $25 \times 25$ , 371  
372 a stark stratification emerges: models with explicit 373  
374 reasoning mechanisms demonstrate exponential 375  
376 depth growth while standard variants plateau. 377

378 The *DeepSeek-V3.2-Thinking* model dominates 379  
380 across all scales, achieving 74.8 steps at  $15 \times 15$  381  
382 and 111.4 at  $25 \times 25$ —3.08 $\times$  deeper than *GPT- 383*  
384 *o3* (36.2 steps). Crucially, this sustained depth 385  
386 persists despite in larger grids, indicating that the 387  
388 “thinking” mechanism enables systematic state- 389  
390 space exploration even when solutions remain 391  
392 elusive. 393

394 In contrast, standard models like *DeepSeek- 395*  
396 *V3.2 (temperature=0)* and *GPT-5.1* exhibit near- 397  
398 flat performance curves, with depths below 15 399  
400 steps even at  $25 \times 25$ . The temperature-ablated 401  
402 configuration shows critically low depth (2.8 steps 403  
404 at  $25 \times 25$ ), directly correlating with its abrupt 405  
406 success rate collapse beyond  $13 \times 13$  grids. This 407  
408 establishes stochastic exploration as a necessary 409  
410 condition for maintaining cognitive depth in com- 411  
412 plex environments. 413

414 The critical divergence between success rate and 415  
416 reasoning depth reveals a key insight: while non- 417  
418 reasoning models fail both metrics at large scales, 419  
420 the *Thinking* variant’s sustained depth demon- 421  
422 strates “persistent problem-solving” behavior— 423  
424 continuing meaningful exploration even when 425

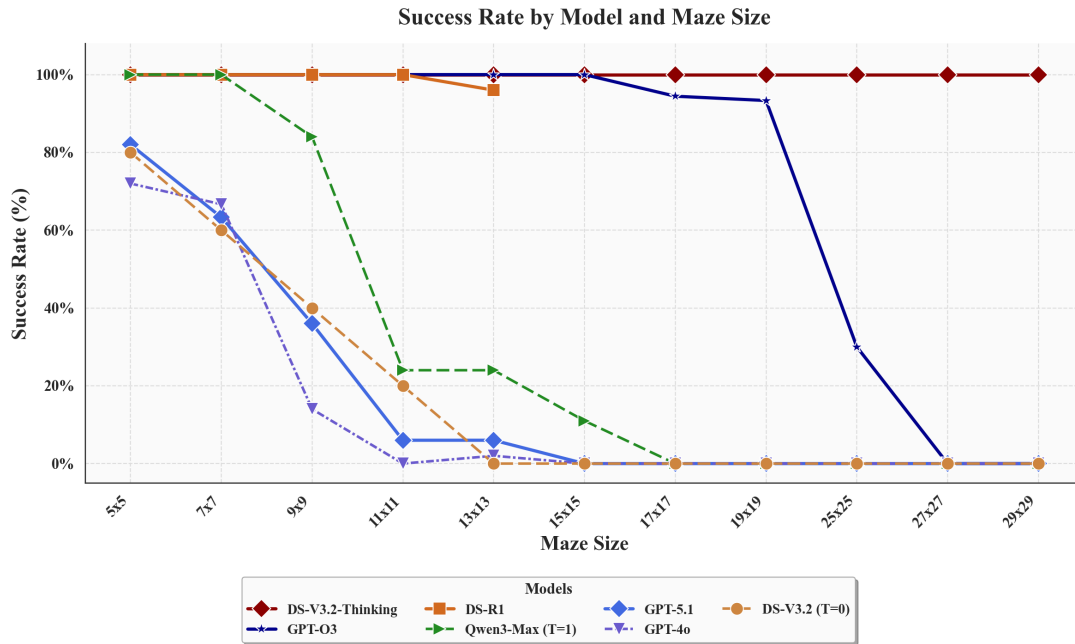


Figure 3: Success rate degradation across maze complexities for different LLM architectures. The plot demonstrates critical performance divergence: reasoning-enabled models (red diamonds) maintain perfect success rates (1.00) through 29×29 grids, while non-reasoning variants (blue circles) exhibit sharp performance cliffs beyond 11×11 grids. Model identifiers: DeepSeek-V3.2-Thinking (red), DeepSeek-V3.2 (temperature=0) (blue), Qwen-3-max-preview (teal), and other variants (light blue).

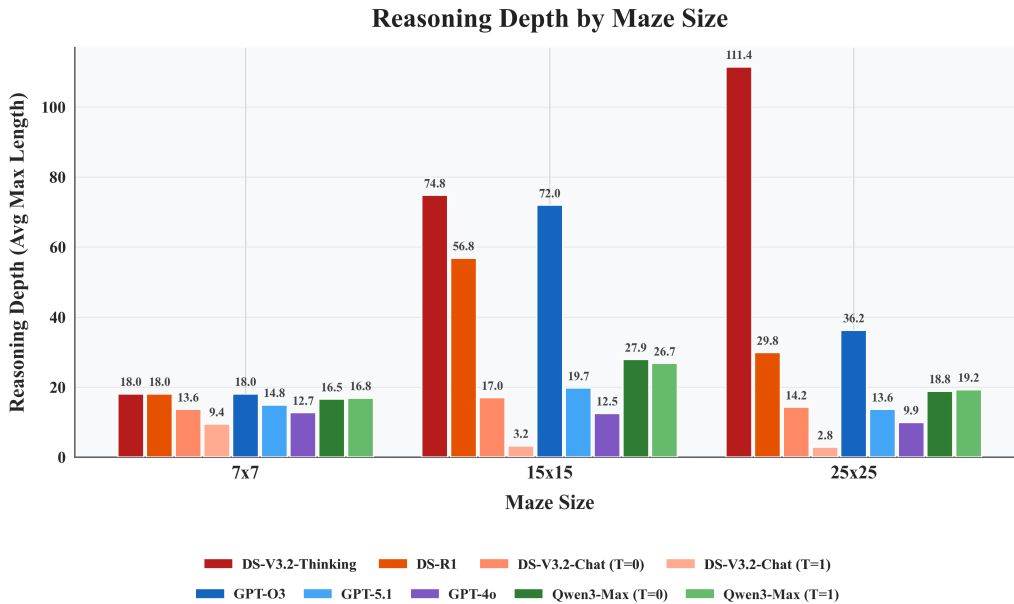


Figure 4: Reasoning depth (in tokenized reasoning steps) across maze complexities. The chart reveals a critical bifurcation: reasoning-enabled models (hatched bars) maintain exponential depth growth (peaking at 111.4 steps for DeepSeek-V3.2-Thinking at 25×25), while non-reasoning variants (solid bars) plateau below 15 steps.

390 success is unattainable within trial constraints. 438  
391 This suggests explicit planning mechanisms fundamen- 439  
392 tally alter cognitive processes by extending 440  
393 the model’s planning horizon rather than merely 441  
394 improving final outcomes. 442

395 These results challenge conventional accuracy- 443  
396 only benchmarks. The depth metric exposes 444  
397 valuable intermediate reasoning capabilities in 445  
398 models that may fail final evaluation—a distinction 446  
399 critical for applications requiring robust partial 447  
400 solutions (e.g., autonomous navigation where partial 448  
401 pathfinding provides actionable data). Future 449  
402 work should investigate whether this depth metric 450  
403 correlates with solution quality when success 451  
404 is achieved, and whether it can inform new 452  
405 training objectives for complex reasoning tasks. 453  
406 The consistent 13×13 grid threshold where non- 454  
407 reasoning models fail both metrics provides a con-  
408 crete benchmark for evaluating spatial reasoning  
409 capabilities.

### 410 4.3.2 Task 2

411 For Task 2, we evaluate models on maze sizes  
412 ranging from  $5 \times 5$  to  $9 \times 9$  with the following  
413 constraints:

- 414 • **Step limits:** Maximum attempts scaled as  $k \times$   
415  $n^2$  where  $n$  is maze size (e.g.,  $1.0x = n^2$  steps,  
416  $2.0x = 2n^2$  steps)
- 417 • **Exploration constraints:** Termination after  
418  $m$  consecutive steps without discovering new  
419 cells (e.g., `noexplore_9` fails after 9 non-  
420 exploring steps)
- 421 • **Visibility radius:** Agent’s perceptual field  
422 size (1-7 cells visible around current position)
- 423 • **Memory window:** Number of recent steps  
424 retained in context (3-7 steps)

425 Fig. 5a shows success rates increase with  
426 visibility radius, with  $7 \times 7$  mazes achieving the  
427 highest success rate (0.5) at visibility radius 7.  $5 \times 5$   
428 and  $9 \times 9$  mazes show lower maximum success rates  
429 (0.36 and 0.2 respectively).

430 Fig. 5b shows success rates decrease as step  
431 limits become stricter (lower  $k$ ) and exploration  
432 patience decreases (lower  $m$ ). For  $9 \times 9$  mazes,  
433 success rate drops from 0.06 (`2.0x_noexplore_9`)  
434 to 0.03 (`1.0x_noexplore_9`), representing a 50%  
435 reduction.

436 Fig. 5c shows  $5 \times 5$  mazes achieve peak success  
437 rate (0.24) at memory window size 7, while  $9 \times 9$

mazes show the lowest success rates across all  
memory window sizes (maximum 0.04 at size 7).

Fig. 5d shows DeepSeek-V3.2-Thinking  
achieves near-perfect success rates (0.99-1.0)  
on  $5 \times 5$  mazes, significantly outperforming  
other models. GPT-o3 demonstrates moderate  
performance (0.75 on  $5 \times 5$ , 0.7 on  $7 \times 7$ , 0.25 on  
 $9 \times 9$ ), while GPT-5.1 and DeepSeek-V3.2 maintain  
consistently low success rates ( $\leq 0.03$ ) across all  
maze sizes. Notably, DeepSeek-V3.2-Thinking  
outperforms all other models by a significant  
margin in all maze dimensions.

The data shows a 3.2× decrease in success rate  
from  $5 \times 5$  to  $9 \times 9$  mazes across most experimental  
conditions. Success rates for  $7 \times 7$  mazes consis-  
tently outperform both smaller and larger mazes in  
visibility and step limit experiments.

## 455 5 Conclusion

### 456 5.1 Conclusion

457 We introduce DEEPMAZE, a minimalist bench-  
458 mark designed to rigorously isolate and quantify  
459 the core spatial reasoning capabilities of large  
460 language models. By leveraging procedurally  
461 generated maze environments with controlled  
462 topology, we eliminate the confounding effects  
463 of prior knowledge exploitation and statistical  
464 shortcuts that plague existing benchmarks. Our  
465 dual-task architecture—comprising *planning* un-  
466 der full observability and *exploration* under partial  
467 observability—forces models to maintain state-  
468 consistent reasoning across sequential actions,  
469 while our novel *Reasoning Depth Metric* explicitly  
470 decouples process quality from outcome success.

471 Experimental validation reveals critical insights  
472 about LLM reasoning capabilities. We identi-  
473 fy a fundamental *reasoning gap*: models with  
474 explicit reasoning mechanisms (e.g., DeepSeek-  
475 V3.2-Thinking) sustain exponential growth in  
476 reasoning depth (up to 111.4 steps) and perfect  
477 success rates even in  $29 \times 29$  mazes, while standard  
478 architectures fail catastrophically beyond  $13 \times 13$   
479 grids despite comparable performance on small  
480 mazes. This bifurcation demonstrates that genuine  
481 multi-step reasoning requires specialized architec-  
482 tural support beyond scaling parameters or training  
483 data. In partial observability settings, we establish  
484 that memory window size and visibility radius are  
485 critical constraints, with success rates dropping  
486 50% when step limits tighten and exploration  
487 patience decreases.

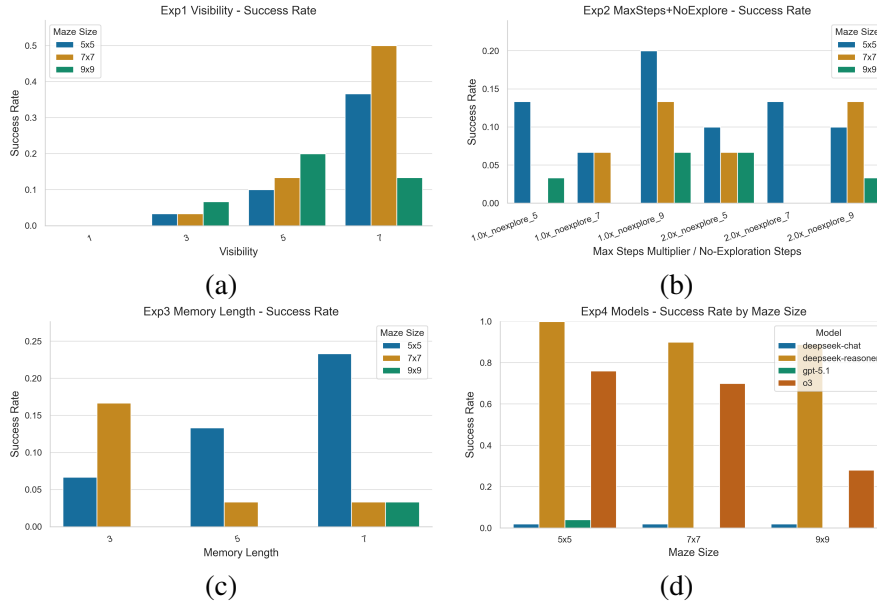


Figure 5: Experimental results across four conditions: (a) Success rate vs. visibility radius (1-7 cells); (b) Success rate under step limits ( $k \times n^2$  steps where  $n$ =maze size) and exploration constraints (failure after  $m$  consecutive non-exploring steps); (c) Success rate vs. memory window size (3-7 steps); (d) Model performance comparison across maze sizes for DeepSeek-V3.2, DeepSeek-V3.2-Thinking, GPT-5.1, and GPT-o3. All experiments evaluate maze-solving success from  $5 \times 5$  to  $9 \times 9$  grids.

DEEPMAZE establishes three foundational principles for reasoning evaluation: (1) **Process over outcome**: Reasoning quality must be measured through intermediate state consistency, not just final answers; (2) **Controlled complexity**: Environmental complexity must be systematically scalable to expose reasoning thresholds; (3) **Diagnostic minimalism**: Core reasoning mechanisms are best isolated by stripping away semantic content and prior-knowledge cues.

These principles position DEEPMAZE not as a domain-specific evaluator, but as a diagnostic layer in the reasoning hierarchy—akin to unit tests in software engineering—that certifies fundamental reasoning competence before deployment in complex environments. Failures here guarantee failures in real-world tasks requiring sustained state-consistent decisions, from supply chain optimization to medical diagnosis.

Future work will extend this framework to: (1) Integrate multimodal observations to bridge the gap between abstract reasoning and embodied cognition; (2) Develop explicit memory architectures that maintain state consistency at scale, informed by our finding that reasoning depth collapses when memory windows fall below critical thresholds; (3) Establish cross-benchmark correlations to validate whether reasoning depth in minimalist environ-

ments predicts performance in complex domain-specific tasks.

As LLMs increasingly operate in safety-critical domains, DEEPMAZE provides a necessary foundation for validating their core reasoning capabilities under controlled conditions—ensuring that apparent competence reflects genuine cognitive depth rather than statistical mirages.

516  
517  
518  
519  
520  
521  
522  
523

## 524 Limitations

525 While DeepMaze provides rigorous diagnostics for  
526 reasoning depth, four boundaries define its scope:

527 (1) **Multimodal grounding:** Our minimal envi-  
528 ronment intentionally excludes visual, auditory, or  
529 tactile inputs, as well as physical dynamics (e.g.,  
530 time constraints, resource limitations). Conse-  
531 quently, DeepMaze *cannot evaluate* reasoning inte-  
532 grated with sensory processing or real-world oper-  
533 ational constraints—capabilities requiring domain-  
534 specific simulators or physical testbeds.

535 (2) **Domain knowledge integration:** The  
536 absence of semantic content (e.g., conceptual  
537 relationships, factual priors) means DeepMaze  
538 isolates *pure structural reasoning*, but cannot as-  
539 sess how domain knowledge interacts with logical  
540 inference—a dimension critical for applications  
541 like scientific reasoning or strategic planning.

542 (3) **Memory mechanism analysis:** Although  
543 sustained reasoning inherently relies on memory,  
544 our framework currently lacks tools to directly  
545 probe *how models utilize memory during rea-*  
546 *soning*. We quantify reasoning depth through  
547 behavioral outcomes but cannot dissect memory  
548 operations (e.g., encoding, retrieval, consolida-  
549 tion) that enable state consistency. This gap  
550 prevents us from establishing causal links between  
551 architectural memory mechanisms and reasoning  
552 performance degradation.

553 (4) **Model coverage:** Evaluations currently  
554 exclude some latest closed-source models (e.g.,  
555 Gemini-series), though our open benchmark en-  
556 ables immediate validation upon API availability.

557 Critically, these limitations reflect *deliberate*  
558 *scoping choices*, not oversights. DeepMaze targets  
559 a specific layer in the reasoning hierarchy (Fig. 1):  
560 “*It is impossible to validate multi-step reasoning*  
561 *in high-noise environments before verifying its*  
562 *existence in controlled settings.*” Thus, results  
563 should be interpreted as **complementary diag-**  
564 **nostics**—not replacements—for domain-specific  
565 evaluations. The limitation regarding memory  
566 mechanisms particularly motivates our future work  
567 on explicit memory architectures (Section X),  
568 where DeepMaze will serve as the foundational  
569 testbed for isolating how different memory designs  
570 impact reasoning depth. This layered valida-  
571 tion approach establishes necessary prerequisites  
572 before introducing confounding variables from  
573 perception, knowledge, or physical constraints.

## References 574

- 575 Maxime Chevalier-Boisvert, Dzmitry Bahdanau, 575  
576 Salem Lahlou, Lucas Willems, Chitwan Saharia, 576  
577 Thien Huu Nguyen, and Yoshua Bengio. 2019. 577  
578 [Babyai: A platform to study the sample efficiency 578](#)  
579 of grounded language learning. *Preprint*, 579  
580 arXiv:1810.08272. 580
- 581 Maxime Chevalier-Boisvert, Bolun Dai, Mark Towers, 581  
582 Rodrigo de Lázcano, Lucas Willems, Salem Lahlou, 582  
583 Suman Pal, Pablo Samuel Castro, and Jordan 583  
584 Terry. 2023. [Minigrad & mineworld: Modular & 584](#)  
585 customizable reinforcement learning environments 585  
586 for goal-oriented tasks. *Preprint*, arXiv:2306.13831. 586
- 587 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, 587  
588 Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias 588  
589 Plappert, Jerry Tworek, Jacob Hilton, Reiichiro 589  
590 Nakano, Christopher Hesse, and John Schulman. 590  
591 2021. [Training verifiers to solve math word 591](#)  
592 problems. *Preprint*, arXiv:2110.14168. 592
- 593 Alan Dao and Dinh Bach Vu. 2025. Alphamaze: En- 593  
594 hancing large language models’ spatial intelligence 594  
595 via grpo. *arXiv preprint arXiv:2502.14669*. 595
- 596 Hafsteinn Einarsson. 2025. Mazeeval: A benchmark 596  
597 for testing sequential decision-making in language 597  
598 models. *arXiv preprint arXiv:2507.20395*. 598
- 599 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul 599  
600 Arora, Steven Basart, Eric Tang, Dawn Song, and 600  
601 Jacob Steinhardt. 2021. [Measuring mathematical 601](#)  
602 problem solving with the math dataset. *Preprint*, 602  
603 arXiv:2103.03874. 603
- 604 Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, 604  
605 Yile Wang, and Yue Zhang. 2020. [Logiqa: A chal- 605](#)  
606 lenge dataset for machine reading comprehension 606  
607 with logical reasoning. *Preprint*, arXiv:2007.08124. 607
- 608 Nora S Newcombe and Thomas F Shipley. 2014. 608  
609 Thinking about spatial thinking: New typology, 609  
610 new assessments. In *Studying visual and spatial 610*  
611 *reasoning for design creativity*, pages 179–192. 611  
612 Springer. 612
- 613 Xavier Puig, Kevin Ra, Marko Boben, Jiaman Li, 613  
614 Tingwu Wang, Sanja Fidler, and Antonio Torralba. 614  
615 2018. [Virtualhome: Simulating household activities 615](#)  
616 via programs. *Preprint*, arXiv:1806.07011. 616
- 617 Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, 617  
618 Yonatan Bisk, Adam Trischler, and Matthew 618  
619 Hausknecht. 2021. [ALFWorld: Aligning Text and 619](#)  
620 [Embodied Environments for Interactive Learning](#). 620  
621 In *Proceedings of the International Conference on 621*  
622 *Learning Representations (ICLR)*. 622
- 623 Peiran Xu, Sudong Wang, Yao Zhu, Jianing Li, and 623  
624 Yunjian Zhang. 2025. [Spatialbench: Benchmarking 624](#)  
625 multimodal large language models for spatial cog- 625  
626 nition. *Preprint*, arXiv:2511.21471. 626

Kaiyu Yang, Aidan M. Swope, Alex Gu, Rahul Chalamala, Peiyang Song, Shixing Yu, Saad Godil, Ryan Prenger, and Anima Anandkumar. 2023. [Leandojo: Theorem proving with retrieval-augmented language models](#). *Preprint*, arXiv:2306.15626.

Weihao Yu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. 2020. [Reclor: A reading comprehension dataset requiring logical reasoning](#). *Preprint*, arXiv:2002.04326.

## A Discussion: DeepMaze as a Diagnostic Layer for Reasoning Capability

Our benchmark intentionally adopts minimalist grid environments to isolate reasoning depth—a design choice that requires explicit contextualization against complex reasoning tasks. Crucially, DeepMaze functions not as a standalone evaluator, but as a controlled diagnostic layer that reveals core reasoning flaws before deployment in complex environments (Fig. 1), analogous to unit tests in software engineering that expose fundamental logic errors prior to system integration. This perspective reframes apparent “toy” constraints as methodological strengths: failures in DeepMaze certify fundamental gaps that persist across diverse reasoning domains.

To formalize this mapping, Table 2 aligns our abstractions with constraints shared across real-world reasoning tasks:

Mechanism	Core Constraint	Real-World Manifestation
State-consistent action sequences	Strict state transition dependency	Sequential decision-making under evolving conditions (e.g., supply chain optimization with dynamic constraints)
Partial observability in exploration	Limited information access	Medical diagnosis with incomplete patient history or financial forecasting with partial market data
Topology complexity scaling	Combinatorial state space growth	Strategic planning in complex systems (e.g., logistics networks with interdependent nodes)

Table 2: DeepMaze isolates reasoning primitives shared across domains. Failures here (e.g., state inconsistency at high complexity) *guarantee* failure in corresponding real-world scenarios, regardless of domain semantics.

Thus, DeepMaze exposes reasoning prerequisites that must be satisfied before deploying models

in resource-intensive applications. Its value lies not in replacing domain-specific evaluation, but in pre-screening models for fundamental reasoning competence—filtering out architectures prone to statistical shortcuts before costly specialized testing. This diagnostic capability becomes increasingly critical as models are deployed in safety-sensitive domains where reasoning failures carry significant consequences.

## B Prompt Templates for MazeBench Tasks

### B.1 Task 1: Global Observation Mode

#### Task 1 Prompt: Global Observation Mode

You are a maze expert and need to navigate yourself to the goal. Current position  $A\{current\_position\}$ , goal  $\{goal\}$ . Coordinate format: All coordinates use  $(y, x)$  format, where  $y$  is the row number and  $x$  is the column number.

You can:

Output path coordinates: e.g.,  $(1,2), (1,3), (2,3)$

Note: Please start outputting from the next position's coordinate, do not output coordinates that duplicate the current position, and the output coordinates must be adjacent.

Hint: First confirm your current position, then you can move to a safe position first, and then plan the next steps.

Maze layout:

$\{maze\_ascii\}$

Symbol legend:

# : Wall

O : Path

S : Start position (you have left)

G : Goal position

A : Current position

? : Unknown area

$\{memory\_info\}$

Output protocol:

Reply with ONLY the path coordinate sequence in the format ``(y1,x1),(y2,x2),...``  
Do not include any explanations, reasoning, or additional text.

670

671

## B.2 Task 2: Partial Observation Mode

### Task 2 Prompt: Partial Observation Mode

You are a maze expert and need to navigate yourself to the goal. You can only see the surrounding {visibility} cells, and your vision cannot pass through walls. You do not know your current coordinates or the goal coordinates.

You can move by specifying directions:

- Direction: up, down, left, right
- Format: direction1,direction2,...
- Each direction moves one step

Constraint 1 -- Determine legal actions first:

Before choosing moves, you must internally determine which directions are immediately blocked by walls or out of bounds. You must NEVER choose a direction that is known to be blocked.

Constraint 2 -- Forbidden actions from failures:

Any direction that resulted in a wall collision in recent steps must be treated as forbidden unless no other legal direction exists.

Constraint 3 -- Exploration requirement:

You must explore new positions. Repeatedly revisiting the same locations will cause the task to be judged as failed.

Visible maze layout (you are at position A):

672

{maze\_ascii}

Symbol legend:

# : Wall  
O : Path  
S : Start position (you have left)  
G : Goal position  
A : Current position  
? : Unknown area

{memory\_info}

Output protocol:

Reply with ONLY the move directions in the format  
``direction1,direction2,...``  
Do not include any explanations, reasoning, or additional text.

673