

What Do LLM Agents Know About Their World? Task2Quiz: A Paradigm for Studying Environment Understanding

Anonymous ACL submission

Abstract

Large language model (LLM) agents have demonstrated remarkable capabilities in complex decision-making and tool-use tasks, yet their ability to generalize across varying environments remains a under-examined concern. Current evaluation paradigms predominantly rely on trajectory-based metrics that measure task success, while failing to assess whether agents possess a grounded, transferable model of the environment. To address this gap, we propose Task-to-Quiz (T2Q), a deterministic and automated evaluation paradigm designed to decouple task execution from world-state understanding. We instantiate this paradigm in T2QBench, a suite comprising 30 environments and 1,967 grounded QA pairs across multiple difficulty levels. Our extensive experiments reveal that task success is often a poor proxy for environment understanding, and that current memory mechanism can not effectively help agents acquire a grounded model of the environment. These findings identify proactive exploration and fine-grained state representation as primary bottlenecks, offering a robust foundation for developing more generalizable autonomous agents.

1 Introduction

Large language model (LLM) agents have recently shown impressive competence in realistic, multi-step decision making and tool use (Yehudai et al., 2025; Luo et al., 2025; Matarazzo and Torlone, 2025; Minaee et al., 2025), enabling progress in web navigation (Wei et al., 2025), software engineering (Zhang et al., 2024; Dong et al., 2025), and operating-system interaction (Hu et al., 2025). While, the issue of generalization is concerned by the community. Agent capabilities that look strong on particular scenarios often fail to transfer reliably when the environment changes or constraints are slightly modified(Liu et al., 2025). This gap raises a basic but under-examined question: do

LLM agents actually acquire a grounded model of the environment, or do they primarily learn specific heuristics that optimize task metrics?

In this paper, we aim at investigating agent’s environment understanding beyond current task success evaluation. As shown in Figure 1, most existing agent benchmarks are trajectory-based. They ask whether an agent reaches the goal, such as success rate or the quality of intermediate actions, measuring “doing”. However, generalization or transferability requires “knowing”. An agent may complete a task while still lacking robust world-state knowledge, such as where objects are, how rooms connect, which direction relationships hold, or what latent object states (e.g., locked/closed) are. Conversely, an agent may acquire correct environment facts yet fail due to planning errors or long-horizon constraints. Therefore, we highlight environment-based evaluation to provide diagnostic “knowing” signals, support controlled interventions (e.g., varying layout complexity or reachability constraints), and enables reproducible grading grounded in known world metadata rather than human annotation or stochastic LLM-as-judge scoring. However, there are three challenges: (i) how to encourage agents to explore beyond single-goal exploitation; (ii) how to quantify environment understanding in a fine-grained, comparable way; and (iii) how to build a fully automated and reproducible pipeline without manual labeling or judge hallucinations (Kalai et al., 2025).

To address the issues, we propose **Task-to-Quiz (T2Q)**, an automated environment-based evaluation paradigm. We instantiate T2Q on TextWorld-style text games, where the full environment metadata (topology, entity placements, and symbolic state relations) is available by construction, making verification precise and convenient. There are two stages. Stage 1 synthesizes coverage-oriented task sets that drive agents to traverse rooms, interact with objects, and reveal stateful

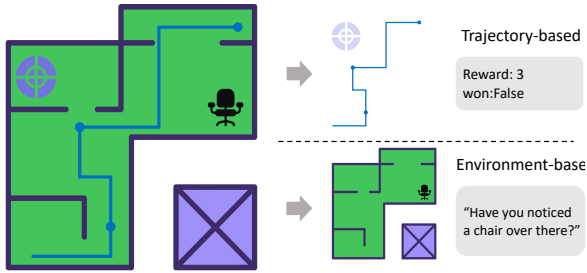


Figure 1: Comparison of trajectory-based and environment-based Evaluation, the trajectory-based evaluation cares about whether and how the final goal is reached, while the environment-based evaluation cares about the entire understanding of the environment.

constraints; task generation is designed to increase exposure to diverse environment elements rather than only optimizing a single end goal. Stage 2 converts the environment and the executed trajectories into a multi-dimensional quiz that probes complementary aspects of world-state knowledge, including localization, connectivity, direction, key-lock matching, and latent properties. Importantly, each question is paired with trajectory-based prerequisites: if an agent’s interactions could not have revealed the necessary evidence (e.g., it never visited a room or never opened a container), the question is marked non-answerable to avoid penalizing agents for information they had no opportunity to observe. Because answers are computed by a verifier grounded in environment metadata and interaction logs, the entire evaluation is deterministic and reproducible. This design separates failures of exploration/exposure from failures of state tracking and reasoning, enabling a fair and diagnostic measure of world-state knowledge.

Using this paradigm, we build T2QBench, a controlled evaluation suite spanning three difficulty levels with 30 environments, 224 coverage-oriented tasks, and 1,967 environment-grounded QA pairs. We evaluate diverse proprietary and open-source models under a unified protocol that jointly measures two complementary dimensions: Task Success Rate (TSR) and Environment Understanding Score (EUS), capturing post-interaction world-state knowledge. Across extensive experiments with different memory mechanisms, we obtain three consistent insights. First, task success is indeed not a reliable proxy for environment understanding. TSR and EUS can diverge as difficulty increases. Second, existing memory systems do not consistently improve environment under-

standing; in many cases, a naive in-context baseline matches or outperforms memory-augmented agents, suggesting that current memory pipelines may lose fine-grained evidence. Third, a low propensity for proactive exploration is a dominant bottleneck. Agents perform relatively better on questions answerable via short-term recall, but struggle on aspects that require actively uncovering latent properties or relations. Together, these results suggest that improving generalization may require not only better planning and retrieval, but also mechanisms that explicitly support world-state acquisition and representation.

Our contributes can be summarized as follows:

- We propose a deterministic environment-based evaluation paradigm **Task-to-Quiz (T2Q)**, that separates “doing” from “knowing”. T2Q contains two stages: 1) coverage-oriented task generation and 2) quiz evaluation, whose answer is dynamically generated based on the agent’s interaction history and the environment metadata.
- We construct a fully automated pipeline and build **T2QBench**, enabling reproducible measurement on different environments, tasks, and difficulty levels.
- We conduct an empirical investigation across models and memory systems that reveals key bottlenecks in agents’ environment understanding ability.

2 Controllable Environment Construction

In this section, we describe our controllable environment construction pipeline. We first present the base environment generation procedure, and then introduce the task coverage planning pipeline.

2.1 Environment Generation

Layout and Connectivity Generation We build the base game environments on top of TextWorld Framework (Côté et al., 2019) with a controllable configuration space. Specifically, we procedurally generate room layouts and connectivity and instantiate an explicit number of interactive objects. As a result, the full environment metadata is known by construction (e.g., connectivity, locations, and object attributes) for future verification and question generation.

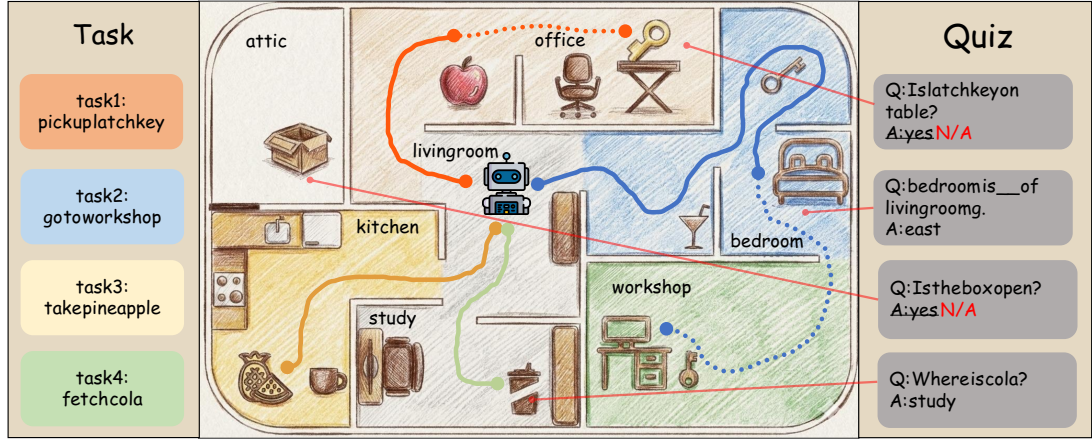


Figure 2: Overview of our T2Q paradigm. It contains coverage-oriented tasks, and environment-based quizzes that are used to evaluate the agent’s environment understanding. The answer of each question is dynamically generated based on the agent’s interaction history and the environment metadata.

Reachability Design The default TextWorld generation pipeline provides limited control over reachability because locks and keys are often sampled randomly, which may lead to weak coupling between access constraints and traversal structure. To make reachability controllable and systematically challenging, we add a post-generation refinement step that injects lock–key dependencies while preserving solvability: we first compute the reachable region from the start, then lock a fixed fraction of doors/containers (0.4 at our default setting) to enforce multi-step, long-horizon planning.

Distractor Placement To further increase reasoning difficulty without altering overall solvability, distractors are added during world construction. These include irrelevant objects and keys that are interactable but useless for task completion.

2.2 Task Coverage Planning

Modeling Task Generation as Weighted Set Cover Problem To comprehensively evaluate an agent’s environment understanding capability, we aim to maximize its opportunities to interact with the environment. Since the exploration stage is task-driven, we design a task set that covers as much of the reachable area and interactive entities as possible while minimizing redundancy among tasks. Therefore, task generation is formulated as a **weighted set cover problem** (Algorithm 1). First, all reachable rooms and interactive entities are extracted from the environment instance \mathcal{G} to form the target universe \mathcal{R} . They are considered as the targets that final task set should cover. For each target in \mathcal{R} , a corresponding target-specific

Algorithm 1 Coverage Task Generation

Require: Environment instance \mathcal{G} , metadata/knowledge \mathcal{K}
Ensure: Greedy task set \mathcal{Q}

- 1: $\mathcal{R} \leftarrow \text{EXTRACTTARGETS}(\mathcal{G})$ \triangleright Universe of targets
- 2: $\mathcal{S} \leftarrow \text{PLANPATHS}(\mathcal{G}, \mathcal{K})$ \triangleright Candidate Trajectories
- 3: $\mathcal{Q} \leftarrow \emptyset$
- 4: $\mathcal{U} \leftarrow \mathcal{R}$ \triangleright Uncovered targets
- 5: **while** $\mathcal{U} \neq \emptyset$ **do**
- 6: $s^* \leftarrow \text{null}, g_{max} \leftarrow 0$
- 7: **for** $s \in \mathcal{S} \setminus \mathcal{Q}$ **do**
- 8: $\Delta \leftarrow \text{COVER}(s) \cap \mathcal{U}$ \triangleright Marginal coverage
- 9: $g \leftarrow \text{GAIN}(\Delta)$ \triangleright Type-weighted gain
- 10: **if** $g > g_{max}$ **then**
- 11: $g_{max} \leftarrow g$
- 12: $s^* \leftarrow s$
- 13: **end if**
- 14: **end for**
- 15: **if** s^* is null **then**
- 16: **break**
- 17: **end if**
- 18: $\mathcal{Q} \leftarrow \mathcal{Q} \cup \{s^*\}$
- 19: $\mathcal{U} \leftarrow \mathcal{U} \setminus \text{COVER}(s^*)$
- 20: **end while**
- 21: **return** $\text{TOTASKS}(\mathcal{Q})$

candidate task is instantiated. To accomplish each task goal, the agent must execute a sequence of actions; we refer to the corresponding reference solution trajectory as a *walkthrough*. Accordingly, a winning walkthrough is derived via TextWorld task analysis function using available metadata \mathcal{K} , producing a candidate task set \mathcal{S} . These generated tasks serve as the candidate pool for the subsequent coverage selection step.

Coverage Signature When executing a task, an agent necessarily visits a subset of rooms and interacts with a subset of entities; we treat these visited and exercised elements as the task’s *coverage signature*. For example, the task “fetch an apple

215	from the kitchen” may require traversing the liv-	the apple in the kitchen?”. This process yields	265
216	ing room and bedroom and opening a chest, yield-	a diverse quiz set that covers environment details	266
217	ing a signature that covers {living room, bedroom,	across multiple question types.	267
218	chest}. Formally, each candidate task $s \in \mathcal{S}$ is as-		
219	sociated with a coverage signature $\text{COVER}(s) \subseteq$	Prerequisite Checkpoints In order to construct	268
220	\mathcal{R} , which includes both the rooms traversed along	trajectory-conditioned, dynamic answers, prereq-	269
221	the path and the interaction-relevant entities exer-	uisite checkpoints are generated together with the	270
222	cised by the task. Therefore, selecting a compact	questions. These checkpoints specify minimal in-	271
223	subset of tasks whose union covers \mathcal{R} naturally in-	teraction evidence required for an agent to be ex-	272
224	duces a (weighted) set cover formulation.	pected to know the answer. For example, a loca-	273
		tion question such as “the apple is in the kitchen”	274
225	Greedy Selection To improve data efficiency,	is only marked answerable if the agent has reached	275
226	we aim to minimize redundancy across tasks and	the kitchen.	276
227	keep the selected task set as compact as possible.		
228	We therefore adopt a greedy selection strategy that,	Dynamic Ground Truth Generation Answers	277
229	at each step, chooses the task that contributes the	are assigned by combining static environment	278
230	largest marginal coverage over the yet-uncovered	facts with the agent’s interaction history. As Fig-	279
231	targets. A greedy weighted selection procedure is	ure 3 shows, if all prerequisite checkpoints are	280
232	adopted in practice. At each iteration, the remain-	satisfied, the agent is considered to have had suf-	281
233	ing uncovered set \mathcal{U} is maintained, the marginal	ficient opportunity to acquire the relevant knowl-	282
234	coverage is computed as $\Delta = \text{COVER}(s) \cap \mathcal{U}$.	edge, and the reference answer is used as the	283
235	Then, the task maximizing a type-weighted gain	ground truth, which is regarded as “ should have	284
236	$\text{GAIN}(\Delta)$ is selected. The gain function assigns	known ” by the agent. Otherwise, the question is	285
237	larger weights to interactions (w_i) and objects	labeled as non-answerable for that trajectory to	286
238	(w_o) than rooms (w_r) in a type-wise manner, pri-	avoid penalizing agents for information they could	287
239	oritizing tasks that exercise complex behaviors over	not have observed. This hybrid of static answers	288
240	pure navigation. The procedure iterates until $\mathcal{U} =$	and trajectory-conditioned answerability reached	289
241	\emptyset , yielding a minimal (approximate) task set Q_{\min}	a balance between fairness and dynamic, behavior-	290
242	with full target coverage and substantially reduced	aware evaluation.	291
243	redundancy.		
244	2.3 Quiz Generation	3 Evaluation	292
245	Beyond coverage-oriented tasks, a complemen-	3.1 Dataset	293
246	tary quiz set is constructed to probe fine-	T2QBench We construct T2QBench via our	294
247	grained grounded understanding of the environ-	pipeline and stratify environments into three dif-	295
248	ment. Question–answer pairs are generated from	difficulty levels based on the number of rooms and	296
249	multiple perspectives.	interactive objects. For each difficulty, we gener-	297
		ate 10 distinct environments. Each environment is	298
250	Convert Environment Facts to Quiz Questions	paired with a coverage-oriented task set and a quiz	299
251	To automatically construct a comprehensive quiz	set. In total, the benchmark comprises 30 environ-	300
252	suite that reflects an agent’s understanding of and	ments, 224 tasks, and 1,967 questions.	301
253	responsiveness to the environment, we generate		
254	quiz questions using deterministic rules grounded	Quiz Taxonomy We categorize quiz questions	302
255	in the environment specification. Concretely, each	into five types to probe complementary facets of	303
256	environment instance is defined by a set of sym-	grounded understanding:	304
257	bolic <i>facts</i> that encode fine-grained details such		
258	as entity locations, object attributes, room orien-	• Location (Loc.): asks where an entity is situ-	305
259	tations, and connectivity (see Section A). We ex-	ated (e.g., the room/container of an object),	306
260	tract a fact list that covers all such environment	which is often answerable given direct obser-	307
261	details and convert each fact into one or more	vation evidence.	308
262	question–answer pairs. For example, the fact		
263	$\text{in}(\text{apple}, \text{kitchen})$ can be transformed into	• Connectivity (Conn.): queries the existence	309
264	questions like “Where is the apple?” and “Is	of traversable links between locations under	310
		the environment topology.	311

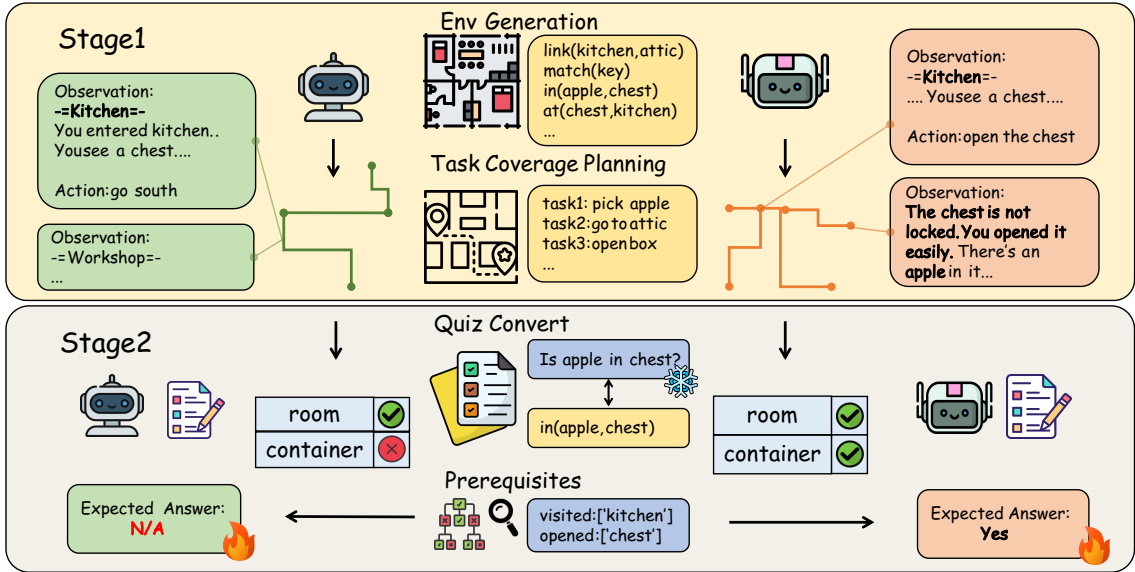


Figure 3: Pipeline of data construction and 2-stage evaluation. Data is generated in deterministic rules during multiple steps, such as room layout, object placement, task coverage planning, and quiz generation. Agent needs to complete a set of tasks first. Then, the agent’s interaction history and the environment metadata are used to generate the quiz set. Finally, the agent is asked to answer the quiz set based on their memory on stage 1.

- **Direction (Dir.):** tests spatial orientation between locations typically demanding cross-turn aggregation and consistent coordinate reasoning.
- **Match (Match.):** tests the agent’s ability to determine which key is the correct one for the lock.
- **Property (Prop.):** concerns latent object states, such as locked/unlocked, open/closed. Agent needs to interact with the object to reveal the property. For example, agent can only reveal a box is locked after trying to open it.

This taxonomy aligns generation with evaluation: each question is paired with trajectory-dependent prerequisites as described in Section 2.3 so that type-wise accuracy reflects reasoning/memory given sufficient exposure rather than missing exploration.

3.2 Formulation of Metrics

Stage 1: Task Success Rate In Stage 1, we evaluate an agent’s performance on the generated task set, reporting TSR. Given an environment \mathcal{E} , we instantiate a task set $\mathcal{T} = \{\tau_1, \tau_2, \dots, \tau_{|\mathcal{T}|}\}$ by the coverage-oriented task generation logic as Algorithm 1. We then evaluate an agent (policy) π by executing each task τ_i in \mathcal{E} . The

resulting interaction trace is a trajectory $\zeta_i = (s_0, a_0, s_1, a_1, \dots, s_{T_i})$, and we denote the set of trajectories as $\mathcal{Z} = \{\zeta_1, \zeta_2, \dots, \zeta_{|\mathcal{T}|}\}$. We define a task success indicator as $\mathbb{I}_{\text{win}}(\zeta_i, \tau_i) \in \{0, 1\}$, and report the task success rate:

$$\text{TSR} = \frac{1}{|\mathcal{T}|} \sum_{i=1}^{|\mathcal{T}|} \mathbb{I}_{\text{win}}(\zeta_i, \tau_i). \quad (1)$$

where $|\zeta_i|$ is the number of actions in ζ_i , and $\ell^*(\tau_i)$ is the oracle shortest action length to complete τ_i in \mathcal{E} (when defined). For memory-augmented agents, Stage 1 interactions also update an internal memory state; we denote the resulting memory after executing the task set (associated with \mathcal{Z}) as $M_{\mathcal{Z}}$, which is carried into Stage 2.

Stage 2: Environment Understanding Score

In Stage 2, we assess grounded understanding with quizzes spanning orientation, connectivity, object properties, and entity localization; importantly, a subset of questions depends on the agent’s realized trajectory and thus requires interaction-conditioned reference answers. We further instantiate a question set $\mathcal{Q} = \{q_1, q_2, \dots, q_{|\mathcal{Q}|}\}$ to probe the agent’s grounded understanding of \mathcal{E} . Crucially, many questions are *trajectory-dependent*; therefore we construct reference answers by answerability analysis $g(\cdot)$ using both the question

and the realized interaction:

$$a_j^* = g(q_j, \mathcal{E}, \mathcal{Z}), \quad \mathcal{A}^* = \{a_k^*\}_{k=1}^{|\mathcal{Q}|}. \quad (2)$$

Given $(q_j, M_{\mathcal{Z}})$, the agent produces a response with $\hat{a}_j = f_{\pi}(q_j, M_{\mathcal{Z}})$, yielding $\hat{\mathcal{A}} = \{\hat{a}_1, \hat{a}_2, \dots, \hat{a}_{|\mathcal{Q}|}\}$. We compute the Environment Understanding Score (EUS) as the ratio of the number of correct answers to the total number of questions:

$$\text{EUS} = \frac{1}{|\mathcal{Q}|} \sum_{j=1}^{|\mathcal{Q}|} \mathbb{I}(\hat{a}_j = a_j^*). \quad (3)$$

Table 1: Dataset Statistics for T2QBench. For each difficulty level, we generate 10 environments with varying numbers of rooms and objects. Walkthrough is the number of steps in a golden solution to complete a task.

Setting	Easy	Medium	Hard
Rooms	3–5	6–10	16–20
Objects	6–10	14–18	28–32
Avg. Tasks	2.4	5.7	12.3
Avg. Walkthrough	4.17	5.25	5.78
Avg. Questions	34.6	65.6	94.7

4 Experiment

This section describes our experimental setup, including the evaluated models and agent baselines, followed by a discussion of the main results and associated analyses.

4.1 Models and Agent Baselines

We evaluate a strong proprietary model GPT-5.1 (OpenAI, 2025) and representative open-source models (DeepSeekV3.2 (DeepSeek-AI et al., 2025), ChatGLM4.6 (Team et al., 2025), and Qwen3-32B (Yang et al., 2025)). On top of each open-source backbone, we compare multiple agent configurations: a naive *In-context* baseline that retains the full interaction history, and memory-augmented baselines, including Mem0 (Chhikara et al., 2025), LangMem (LangChain, 2024), and A-MEM (Xu et al., 2025) that implement different memory systems.

4.2 Main Result

We report both task completion performance and quiz accuracy on T2QBench across models and memory mechanisms. Quiz results are further

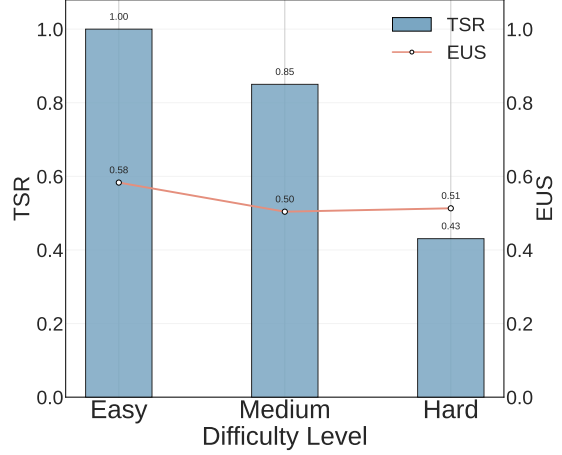


Figure 4: TSR and EUS trends across difficulty levels for GLM-4.6 with the **in-context** method. TSR: task success rate; EUS: environment understanding score.

broken down by question categories (e.g., direction, location, and object properties). Table 2 reveals consistent patterns across models, memory mechanisms, and question types. GPT-5.1 and DeepSeekV3.2 tend to achieve stronger task completion, whereas Qwen3-32B yields the best overall quiz-based environment understanding. In contrast, Mem0 and A-MEM often lag behind on both TSR and EUS, and the simple in-context baseline remains highly competitive. Looking into quiz categories, agents typically perform best on localization questions but struggle the most with object-property queries.

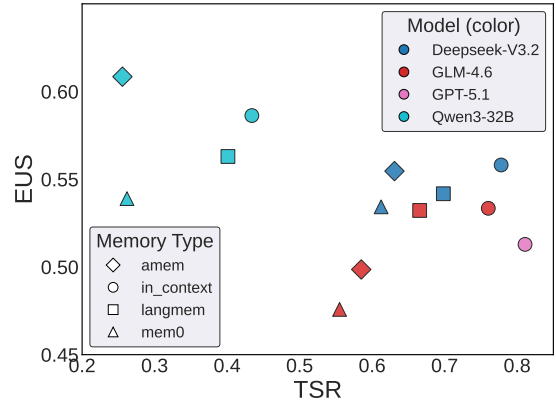


Figure 5: TSR and EUS Comparison: X-axis is the TSR score, Y-axis is the EUS score. The color represents the model, the shape represents the memory system.

4.3 Analysis

In this section, we analyze the main results and figure out the answers to the research questions.

Model	Method	Task Score	Environment Understanding Score					
			Loc.	Conn.	Dir.	Match.	Prop.	Tot.
GPT-5.1		68.75	62.52	51.66	45.68	43.53	35.93	50.23
GLM-4.6	In-context	61.16	62.85	56.40	53.33	38.24	35.65	52.41
	Mem0	39.29	55.32	52.84	50.37	24.71	28.41	46.21
	LangMem	51.34	60.23	59.72	57.04	33.53	30.08	51.65
	A-MEM	42.86	56.46	58.53	52.59	27.65	27.30	48.30
DeepseekV3.2	In-context	63.84	68.41	59.48	53.09	49.41	27.58	54.25
	Mem0	43.30	61.87	60.66	53.58	44.71	28.97	52.41
	LangMem	54.46	66.94	57.35	50.37	44.12	27.86	52.36
	A-MEM	45.54	<u>66.94</u>	59.00	56.05	<u>46.47</u>	30.36	54.55
Qwen3-32B	In-context	34.38	63.38	69.87	58.90	44.30	43.65	58.55
	Mem0	18.30	56.63	66.82	58.52	36.47	42.90	54.96
	LangMem	27.68	61.37	67.77	57.28	37.06	44.01	56.63
	A-MEM	18.30	61.05	71.33	66.67	44.12	46.52	60.29

Table 2: Main experiment results. Best and second-best values in each column are highlighted in **bold** and underlined, respectively, across all model-memory system combinations. Column abbreviations: Loc. (Location), Conn. (Connectivity), Dir. (Direction), Match (Matching), Prop. (Properties), Tot. (Overall).

RQ1: Different Trends of TSR and EUS Across Difficulty Levels To figure out the relationship between TSR and EUS, we evaluate the TSR and EUS trends across different difficulty levels. We observe that **TSR** and **EUS** exhibit markedly different behaviors as environment difficulty increases. Figure 4 illustrates this effect for GLM-4.6 under an in-context baseline: TSR saturates at 100% on easy environments but degrades substantially with increasing difficulty. In contrast, EUS remains comparatively stable across difficulty levels, suggesting that task success is not a reliable proxy for environment understanding.

Takeaway 4.1 for RQ1

Doing is different from knowing. Task success primarily tracks trajectory difficulty (e.g., longer solutions and tighter constraints), whereas EUS captures a more structural form of environment modeling grounded in interaction. Consequently, evaluating agents solely by completion metrics cannot fully reflect environment understanding.

RQ2: Comparison between In-context Method and Memory Systems We then examine the effect of memory systems on environment understanding. We compare a naive in-context baseline with several representative memory-augmented variants, and summarize TSR and EUS across models and memory systems in Figure 5. Surprisingly, the in-context approach often attains the strongest overall performance in both task success

and environment understanding, outperforming methods equipped with explicit memory systems. This suggests that current memory systems do not reliably help agents form a more complete understanding of exploratory environments. A plausible explanation is that, during memory construction and retrieval, high-level abstractions or summaries discard critical fine-grained evidence, resulting in worse downstream reasoning than simply retaining the full interaction context. More broadly, limited environment mastery not only lowers EUS but also harms TSR, indicating that a comprehensive understanding of the environment is also important for tasks.

Takeaway 4.2 for RQ2

Current memory systems helps organization less than it hurts fidelity. Existing memory systems tend to lose environment-specific details while still failing to distill them into structural abstractions such as layout. This gap may stem from event-centric memory systems, highlighting the need for memory mechanisms tailored to grounded world-state representation.

RQ3: Comparison between Different Question Types To further investigate what hinders the agent’s environment understanding, we analyze the performance of different question types in Figure 6. Models achieve relatively high accuracy on *location* questions, which are often directly answerable from a single recent observation. Performance drops for questions like orientation requir-

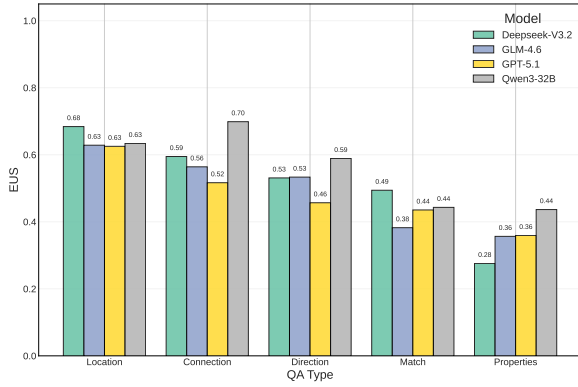


Figure 6: EUS Comparison Across Question Types: Agents have different performance on different question types.

ing multi-turn retrieval and reasoning across interactions. The largest degradation appears on questions that require *proactive* interaction to reveal latent properties or relations. Overall, these patterns suggest that insufficient exploration incentivizes is the dominant bottleneck than retrieval from memory.

Takeaway 4.3 for RQ3

Exploration is the dominant bottleneck than retrieval. Agents are typically optimized to complete a given task with minimal interactions, which incentivizes “efficient” goal-directed behavior. As a result, agents often fail to uncover environment details for deeper understanding. This suggests that the insufficient exploration incentives are the primary bottleneck than retrieval from memory.

5 Related Work

5.1 LLM Agentic Benchmarks

Agentic benchmarks increasingly evaluate LLM agents in *executable* environments that support multi-step interaction and tool use, such as realistic web navigation (Zhou et al., 2023; He et al., 2024) and repository-level software engineering tasks (Jimenez et al., 2023). More broadly, recent suites cover apps (Trivedi et al., 2024), OS/GUI-like scenarios (Agashe et al., 2024; Xie et al., 2024), and multi-environment or generalist evaluation (Liu et al., 2024; Mialon et al., 2023). Despite their achievements, many benchmarks remain task-centric and costly to scale, motivating more controllable and fine-grained evaluation be-

yond final success.

5.2 Game-based Evaluation

Games provide dynamic, multi-skill environments that stress long-horizon planning, exploration, and state tracking. Recent benchmarks include text/conversational games (Qiao et al., 2023; Costarelli et al., 2024), live computer games (Hu et al., 2024; Li et al., 2025), and multimodal games (Paglieri et al., 2024). Strategic and multi-agent settings further evaluate LLM decision-making (Wu et al., 2023; Duan et al., 2024; Chen et al., 2024) and broader gameplay ability (Huang et al., 2024; Chalamalasetti et al., 2023). Text-based adventure games are particularly attractive for grounded evaluation due to explicit state transitions and interpretable actions. Frameworks such as TextWorld (Côté et al., 2019) and interactive fiction resources (Hausknecht et al., 2020; Phan et al., 2025) become an alternative to mirror real-world challenges. In contrast to these evaluation works, we leverage text games as a framework to conduct the entire comprehensive evaluation of the agent’s environment understanding.

6 Conclusion

In this paper, we introduce **Task2Quiz(T2Q)**, a new paradigm and automatically built **T2QBench**, a benchmark spanning three difficulty levels with 30 environments, 224 coverage-oriented tasks, and 1,967 environment-grounded QA pairs. We design a unified two-stage evaluation with Task Success Rate (TSR) and Environment Understanding Score (EUS). We present a comprehensive empirical study on **T2QBench**, and reveal **three key insights**: (i) There is a gap between ‘knowing’ and ‘doing’, task success rate can not reflect what agents *know* about the environment, as TSR degrades with difficulty while EUS remains comparatively stable; (ii) recent memory systems do not achieve better EUS than the naive in-context method, suggesting they are not yet effective at extracting and structuring environment-relevant information from interaction traces; and (iii) It’s the lack of propensity for proactive exploration what hinders agents to form comprehensive world models. Beyond trajectory-level, our work offers a principled route to evaluate environment modeling capability of agents, providing a diagnostic foundation for developing more robust agents with generalized abilities.

529 Limitations

530 Despite our comprehensive analysis of agents' en-
531 vironment understanding, our work has several
532 limitations.

533 First, due to computational constraints and the
534 cost of API-based evaluation, we evaluate only a
535 subset of representative open-source and closed-
536 source models; extending coverage to a broader
537 range of agents can improve the reliability of our
538 findings.

539 Second, our benchmark is built on TextWorld
540 framework (Côté et al., 2019), whose game me-
541 chanics and maps are relatively simple and purely
542 text-based compared with commercial game such
543 as Minecraft. Future extensions could incorporate
544 more complex dynamics, larger maps, and addi-
545 tional modalities, such as visual or audio. There
546 are two possible ways to do this, both challeng-
547 ing: (1) adapting a commercial game as the devel-
548 opment environment, where access to copyright-
549 protected source code is often restricted; and (2)
550 building a new framework from scratch using a
551 more powerful engine such as Unity (Unity Tech-
552 nologies, 2022) or Unreal Engine (Epic Games,
553 2023). The latter would require substantial engi-
554 neering effort comparable to that of a game studio.
555 Although our game framework is simple, it is suf-
556 ficient to support the key insights of our work.

557 Finally, the world-analysis algorithm for
558 coverage-oriented sub-tasks relies on TextWorld
559 framework. It is implemented by a dependency
560 tree structure and exhaustive state-space search
561 (e.g., graph traversal over world states), which
562 can become computationally expensive as the
563 number of rooms and interactive objects increases
564 for more complex games. Improving the effi-
565 ciency and scalability of this algorithm remains a
566 promising avenue for future research.

567 References

568 Saaket Agashe, Jiuzhou Han, Shuyu Gan, Jiachen
569 Yang, Ang Li, and Xin Eric Wang. 2024. *Agent s:*
570 *An open agentic framework that uses computers like*
571 *a human*. *Preprint*, arXiv:2410.08164.

572 Greg Brockman, Vicki Cheung, Ludwig Pettersson,
573 Jonas Schneider, John Schulman, Jie Tang, and Wo-
574 jciech Zaremba. 2016. Openai gym. *arXiv preprint*
575 *arXiv:1606.01540*.

576 Kranti Chalamalasetti, Jana Götze, Sherzod Haki-
577 mov, Brielen Madureira, Philipp Sadler, and David
578 Schlangen. 2023. clembench: Using game play to

579 evaluate chat-optimized language models as conver-
580 sational agents. In *Proceedings of the Conference on*
581 *Empirical Methods in Natural Language Processing*
582 *(EMNLP)*.

583 Junzhe Chen, Xuming Hu, Shuodi Liu, Shiyu Huang,
584 Wei-Wei Tu, Zhaofeng He, and Lijie Wen. 2024.
585 Llmarena: Assessing capabilities of large language
586 models in dynamic multi-agent environments. In
587 *Proceedings of the Annual Meeting of the Associa-*
588 *tion for Computational Linguistics (ACL)*.

589 Prateek Chhikara, Dev Khant, Saket Aryan, Taranjeet
590 Singh, and Deshraj Yadav. 2025. *Mem0: Building*
591 *production-ready AI agents with scalable long-term*
592 *memory*. *CoRR*, abs/2504.19413.

593 Anthony Costarelli, Mat Allen, Roman Hauksson,
594 Grace Sodunke, Suhas Hariharan, Carlson Cheng,
595 Wenjie Li, Joshua Clymer, and Arjun Yadav. 2024.
596 *Gamebench: Evaluating strategic reasoning abilities*
597 *of llm agents*. *arXiv preprint arXiv:2406.06613*.

598 Marc-Alexandre Côté, Ákos Kádár, Xingdi Yuan, Ben
599 Kybartas, Tavian Barnes, Emery Fine, James Moore,
600 Matthew Hausknecht, Layla El Asri, Mahmoud
601 Adada, and 1 others. 2019. Textworld: A learn-
602 ing environment for text-based games. In *Computer*
603 *Games: 7th Workshop, CGW 2018, Held in Conjun-*
604 *ction with the 27th International Conference on Arti-*
605 *ficial Intelligence, IJCAI 2018, Stockholm, Sweden,*
606 *July 13, 2018, Revised Selected Papers*. Springer.

607 DeepSeek-AI, Aixin Liu, Aoxue Mei, Bangcai Lin,
608 Bing Xue, Bingxuan Wang, Bingzheng Xu, Bochao
609 Wu, Bowei Zhang, Chaofan Lin, Chen Dong,
610 Chengda Lu, Chenggang Zhao, Chengqi Deng,
611 Chenhao Xu, Chong Ruan, Damai Dai, Daya Guo,
612 Dejian Yang, and 245 others. 2025. *Deepseek-v3.2:*
613 *Pushing the frontier of open large language models.*
614 *Preprint*, arXiv:2512.02556.

615 Yihong Dong, Xue Jiang, Jiaru Qian, Tian Wang, Kechi
616 Zhang, Zhi Jin, and Ge Li. 2025. *A survey on*
617 *code generation with llm-based agents*. *Preprint*,
618 arXiv:2508.00083.

619 Jinhao Duan, Renming Zhang, James Diffenderfer,
620 Bhavya Kailkhura, Lichao Sun, Elias Stengel-Eskin,
621 Mohit Bansal, Tianlong Chen, and Kaidi Xu. 2024.
622 *Gtbench: Uncovering the strategic reasoning capa-*
623 *bilities of llms via game-theoretic evaluations*. In
624 *Advances in Neural Information Processing Sys-*
625 *tems*.

626 Epic Games. 2023. Unreal engine 5. <https://www.unrealengine.com>. Version 5.3, accessed 2025-01.

629 Matthew J. Hausknecht, Prithviraj Ammanabrolu,
630 Marc-Alexandre Côté, and Xingdi Yuan. 2020. *In-*
631 *teractive fiction games: A colossal adventure*. In
632 *The Thirty-Fourth AAAI Conference on Artificial In-*
633 *telligence, AAAI 2020, The Thirty-Second Inno-*
634 *vative Applications of Artificial Intelligence Confer-*
635 *ence, IAAI 2020, The Tenth AAAI Symposium on Ed-*
636 *ucational Advances in Artificial Intelligence, EAAI*

637	2020, New York, NY, USA, February 7-12, 2020,	Chenwu Liu, Jingyang Yuan, Shichang Zhang, and	692
638	pages 7903–7910. AAAI Press.	7 others. 2025. Large language model agent: A survey on methodology, applications and challenges.	693
639	Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu,	<i>Preprint</i> , arXiv:2503.21460.	694
640	Yong Dai, Hongming Zhang, Zhenzhong Lan, and	Andrea Matarazzo and Riccardo Torlone. 2025. A survey on large language models with some insights on their capabilities and limitations.	696
641	Dong Yu. 2024. Webvoyager: Building an end-to-	<i>Preprint</i> , arXiv:2501.04040.	697
642	end web agent with large multimodal models. <i>arXiv</i>	Grégoire Mialon, Clémentine Fourier, Thomas Wolf,	698
643	<i>preprint arXiv:2401.13919.</i>	Yann LeCun, and Thomas Scialom. 2023. Gaia: a	699
644	Lanxiang Hu, Qiyu Li, Anze Xie, Nan Jiang, Ion Sto-	benchmark for general ai assistants. In <i>The Twelfth</i>	700
645	ica, Haojian Jin, and Hao Zhang. 2024. Gamearena:	<i>International Conference on Learning Representations.</i>	701
646	Evaluating llm reasoning through live computer	Shervin Minaee, Tomas Mikolov, Narjes Nikzad,	702
647	games. <i>arXiv preprint arXiv:2412.06394.</i>	Meysam Chenaghlu, Richard Socher, Xavier Ama-	703
648	Xueyu Hu, Tao Xiong, Biao Yi, Zishu Wei, Ruixuan	trian, and Jianfeng Gao. 2025. Large language models: A survey.	704
649	Xiao, Yurun Chen, Jiasheng Ye, Meiling Tao, Xi-	<i>Preprint</i> , arXiv:2402.06196.	705
650	angxin Zhou, Ziyu Zhao, Yuhuai Li, Shengze Xu,	OpenAI. 2025. GPT-5.1. OpenAI API documentation.	706
651	Shenzhi Wang, Xinchun Xu, Shuofei Qiao, Zhaokai	Accessed via OpenAI API.	707
652	Wang, Kun Kuang, Tiejong Zeng, Liang Wang, and	Davide Paglieri, Bartłomiej Cupiał, Samuel Coward,	708
653	10 others. 2025. Os agents: A survey on mllm-based agents for general computing devices use.	Ulyana Piterbarg, Maciej Wolczyk, Akbir Khan,	709
654	<i>Preprint</i> , arXiv:2508.04482.	Eduardo Pignatelli, Łukasz Kuciński, Lerrel Pinto, Rob	710
655	Jen-tse Huang, Eric John Li, Man Ho Lam, Tian	Fergus, and 1 others. 2024. Balrog: Benchmarking	711
656	Liang, Wenxuan Wang, Youliang Yuan, Wenxiang	agentic llm and vlm reasoning on games. <i>arXiv</i>	712
657	Jiao, Xing Wang, Zhaopeng Tu, and Michael R Lyu.	<i>preprint arXiv:2411.13543.</i>	713
658	2024. How far are we on the decision-making of	Long Phan, Mantas Mazeika, Andy Zou, and Dan	714
659	llms? evaluating llms’ gaming ability in multi-agent	Hendrycks. 2025. Textquests: How good are llms at text-based video games?	715
660	environments. <i>arXiv preprint arXiv:2403.11807.</i>	<i>Preprint</i> , arXiv:2507.23701.	716
661	Carlos E Jimenez, John Yang, Alexander Wettig,	Dan Qiao, Chenfei Wu, Yaobo Liang, Juntao Li,	717
662	Shunyu Yao, Kexin Pei, Ofir Press, and Karthik	and Nan Duan. 2023. Gameeval: Evaluating	718
663	Narasimhan. 2023. Swe-bench: Can language mod-	llms on conversational games. <i>arXiv preprint</i>	719
664	els resolve real-world github issues? <i>arXiv preprint</i>	<i>arXiv:2308.10032.</i>	720
665	<i>arXiv:2310.06770.</i>	5 Team, Aohan Zeng, Xin Lv, Qinkai Zheng, Zhenyu	721
666	Adam Tauman Kalai, Ofir Nachum, Santosh S. Vem-	Hou, Bin Chen, Chengxing Xie, Cunxiang Wang,	722
667	pala, and Edwin Zhang. 2025. Why language models hallucinate.	Da Yin, Hao Zeng, Jiajie Zhang, Kedong Wang,	723
668	<i>Preprint</i> , arXiv:2509.04664.	Lucen Zhong, Mingdao Liu, Rui Lu, Shulin Cao,	724
669	LangChain. 2024. Langmem: Long-term memory for llm agents.	Xiaohan Zhang, Xuancheng Huang, Yao Wei, and	725
670	Accessed: 2025-12-28.	152 others. 2025. Glm-4.5: Agentic, reasoning, and coding (arc) foundation models.	726
671	Zongyuan Li, Yanan Ni, Runnan Qi, Lumin Jiang,	<i>Preprint</i> , arXiv:2508.06471.	727
672	Chang Lu, Xiaojie Xu, Xiangbei Liu, Pengfei Li,	Harsh Trivedi, Tushar Khot, Mareike Hartmann,	728
673	Yunzheng Guo, Zhe Ma, Huanyu Li, Hui Wu, Xian	Ruskin Manku, Vinty Dong, Edward Li, Shashank	729
674	Guo, Kuihua Huang, and Xuebo Zhang. 2025. Llm-pysc2: Starcraft ii learning environment for large language models.	Gupta, Ashish Sabharwal, and Niranjan Balasubra-	730
675	<i>Preprint</i> , arXiv:2411.05348.	manian. 2024. Appworld: A controllable world of	731
676	Weiwen Liu, Jiarui Qin, Xu Huang, Xingshan Zeng,	apps and people for benchmarking interactive cod-	732
677	Yunjia Xi, Jianghao Lin, Chuhan Wu, Yasheng	ing agents. <i>arXiv preprint arXiv:2407.18901.</i>	733
678	Wang, Lifeng Shang, Ruiming Tang, Defu Lian,	Unity Technologies. 2022. Unity real-time devel-	734
679	Yong Yu, and Weinan Zhang. 2025. The real barrier to llm agent usability is agentic roi.	opment platform. https://unity.com . Version	735
680	<i>Preprint</i> , arXiv:2505.17767.	2022.3 LTS, accessed 2025-01.	736
681	Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu	Zhepei Wei, Wenlin Yao, Yao Liu, Weizhi Zhang, Qin	737
682	Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen	Lu, Liang Qiu, Changlong Yu, Puyang Xu, Chao	738
683	Men, Kejuan Yang, and 1 others. 2024. Agentbench:	Zhang, Bing Yin, Hyokun Yun, and Lihong Li.	739
684	Evaluating llms as agents. In <i>ICLR.</i>	2025. Webagent-r1: Training web agents via end-to-end multi-turn reinforcement learning.	740
685	Junyu Luo, Weizhi Zhang, Ye Yuan, Yusheng Zhao,	<i>Preprint</i> , arXiv:2505.16421.	741
686	Junwei Yang, Yiyang Gu, Bohan Wu, Binqi Chen,		742
687	Ziyue Qiao, Qingqing Long, Rongcheng Tu, Xiao		743
688	Luo, Wei Ju, Zhiping Xiao, Yifan Wang, Meng Xiao,		744
689			745
690			746
691			747

748	Yue Wu, Xuan Tang, Tom Mitchell, and Yuanzhi Li.	a controllable underlying graph. The player start	800
749	2023. Smartplay: A benchmark for llms as intelli-	location is also deterministically assigned. Since	801
750	gent agents. In <i>International Conference on Learn-</i>	the topology is programmatically constructed, the	802
751	<i>ing Representations</i> .	ground-truth map (rooms and connectivity) is al-	803
752	Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan	ways available for verification and question gener-	804
753	Li, Siheng Zhao, Ruisheng Cao, Toh J Hua, Zhoujun	ation.	805
754	Cheng, Dongchan Shin, Fangyu Lei, and 1 others.	Entities (object control). TextWorld provides	806
755	2024. Osworld: Benchmarking multimodal agents	typed entities that can be composed to form richer	807
756	for open-ended tasks in real computer environments.	worlds. Rooms and fixed supports define where	808
757	<i>Advances in Neural Information Processing Systems</i> ,	objects can appear, while doors and containers	809
758	37:52040–52094.	introduce stateful constraints (open/closed/locked)	810
759	Wujiang Xu, Zujie Liang, Kai Mei, Hang Gao, Juntao	that shape exploration and interaction. Each entity	811
760	Tan, and Yongfeng Zhang. 2025. A-MEM: agentic	can be named and optionally given a description	812
761	memory for LLM agents . <i>CoRR</i> , abs/2502.12110.	used by examine.	813
762	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,	Facts and relations (state control). Beyond en-	814
763	Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao,	tity placement, the environment dynamics and	815
764	Chengen Huang, Chenxu Lv, Chujie Zheng, Day-	many evaluable properties are encoded as sym-	816
765	iheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao	bolic predicates added via <code>add_fact</code> . This makes	817
766	Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41	the world fully controllable at the level of discrete	818
767	others. 2025. Qwen3 technical report . <i>Preprint</i> ,	state transitions (e.g., whether a door is locked)	819
768	arXiv:2505.09388.	and relational structure (e.g., which key matches	820
769	Asaf Yehudai, Lilach Eden, Alan Li, Guy Uziel, Yilun	which lock), and provides a clean source of envi-	821
770	Zhao, Roy Bar-Haim, Arman Cohan, and Michal	ronment metadata for our verifier.	822
771	Shmueli-Scheuer. 2025. Survey on evaluation of	Supported entity types. We use the standard	823
772	llm-based agents . <i>Preprint</i> , arXiv:2503.16416.	types provided by GameMaker:	824
773	Kechi Zhang, Jia Li, Ge Li, Xianjie Shi, and Zhi	• r: room; d: door; c: container; s: supporter	825
774	Jin. 2024. Codeagent: Enhancing code gener-	• o: portable object; k: key; f: food	826
775	ation with tool-integrated agent systems for real-	• oven, stove: specialized heat-source con-	827
776	world repo-level coding challenges . <i>Preprint</i> ,	tainer/supporter	828
777	arXiv:2401.07339.	Supported predicates (facts). We mainly rely	829
778	Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou,	on the following predicates for controllable state	830
779	Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue	and relations:	831
780	Ou, Yonatan Bisk, Daniel Fried, and 1 others.	• <code>match(key, door/container)</code> : key–lock	832
781	2023. Webarena: A realistic web environment	compatibility	833
782	for building autonomous agents. <i>arXiv preprint</i>	• <code>open(door/container)</code> ,	834
783	<i>arXiv:2307.13854</i> .	<code>closed(door/container)</code> ,	835
784	A The detail of Constructing an	<code>locked(door/container)</code> : access state	836
785	environment	• <code>edible(food)</code> : consumability (optional)	837
786	We build each text-game environment with	How to add new structure (examples). All	838
787	TextWorld’s GameMaker API, which exposes the	controllable components are added deterministi-	839
788	full world state by construction and allows us	cally in code, e.g., (i) create rooms and con-	840
789	to precisely control layout, entities, and stateful	nect them; (ii) create entities and place them in	841
790	relations. Concretely, an environment is speci-	a room/supporter/inventory; and (iii) attach facts	842
791	fied through three steps: (i) defining the world	to specify states and relations. For instance, a	843
792	topology (rooms and connectivity); (ii) instantiat-		
793	ing entities (doors, containers, objects) and plac-		
794	ing them; and (iii) declaring symbolic facts (e.g.,		
795	locked/closed/match) that govern dynamics and		
796	answerable knowledge.		
797	World topology (layout control). We explicitly		
798	create rooms and connect them via directional		
799	edges (e.g., <code>roomA.east</code> to <code>roomB.west</code>), yielding		

locked door and its matching key can be created by adding `locked(door)` and `match(key, door)`, while containers can be initialized with `closed(container)` before gameplay.

Quests and distractors (task control). Given a constructed world, we can (a) record a specific quest by logging an intended command sequence, or (b) automatically sample multiple random quests of bounded length. We can also inject distractor objects to increase interaction complexity without changing the underlying ground-truth world state. Together, these controls allow us to scale environment instances while keeping their structure and metadata fully known and reproducible.

Minimal API examples. Below we provide short code snippets illustrating how we control topology, entities, and symbolic facts using TextWorld’s GameMaker.

Example 1: Define rooms and connectivity (topology control)

```
M = GameMaker()
roomA = M.new_room("Room A")
roomB = M.new_room("Room B")
M.connect(roomA.east, roomB.west)
# directed ports define the map
M.set_player(roomA)
# fix the start room
```

Example 2: Add entities and place them (entity control)

```
s =M.new(type="s",name="table")
roomA.add(s)
# fixed in the room
```

Example 3: Add a locked door and matching key (fact/relation control)

```
d=M.new_door(corridor, name="door")
M.add_fact("locked", d)
k=M.new(type="k", name="old key")
M.add_fact("match", k, d)
s.add(k)
```

Example 4: Add a container with an explicit initial state

```
c=M.new(type="c", name="fridge")
M.add_fact("closed", c)
# explicit initial status
roomB.add(c)
```

Entity types. We use TextWorld’s typed object system, e.g., `r` (room), `d` (door), `c` (container), `s` (supporter), `o` (portable object), and `k` (key).

Controllable predicates. We encode state and relations as symbolic facts, including `open`, `closed`, `locked`, and `match(key, door/container)`. Because these facts are added by construction, the complete environment metadata (topology, placements, and states) is fully known and can be used for deterministic verification and question generation.

B Example of a game environment

Our generated game is packaged as a gym (Brockman et al., 2016) environment. Since all gym games have only one winning state. In order to apply our task sets on one environment, we actually generate a new game environment with the same structure but different goal. So one env \mathcal{E} with n tasks $\mathcal{T} = \{t_1, t_2, \dots, t_n\}$ is actually organized as n games $\mathcal{G}_1 = (\mathcal{E}, t_1), \mathcal{G}_2 = (\mathcal{E}, t_2), \dots, \mathcal{G}_n = (\mathcal{E}, t_n)$. The game content is shown in Figure 7. The task description is shown as follows:

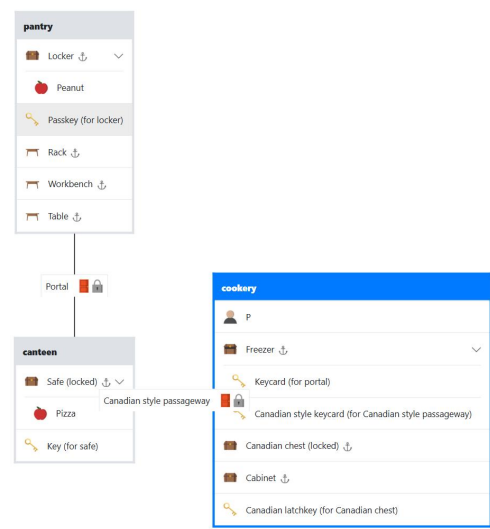


Figure 7: Example of a game environment

Task Description

Get ready to pick stuff up and put it in places, because you've just entered TextWorld! I hope you're ready to go into rooms and interact with objects, Recover the key from the floor of the canteen.

C Prompt Templates

Our prompt templates are shown as follows:

QA Type-specific Guidelines

yes_no: Answer “yes” or “no” only.
which: Answer with the exact choice text from the given choices.
description: Provide a concise, categorical description. Focus on object types/categories rather than specific item names. Do NOT include exits, doors, or room connections. Valid options: [object, container, supporter, food, key].
where: Answer with the location/room name.
what: Answer concisely with the specific information requested.
default: Answer concisely.

Action Prompt Template

Observation:
{obs}

Inventory:
{inventory}

Score: {score}

Return ONLY JSON: {"reason": "...", "command": "<one command>"}

QA Prompt Template

You answer questions about TextWorld gameplay.
Based on the provided context, answer the following question about the TextWorld game initial state.

Context:
{context} (only included if context \neq "")

{question}

Choices:

- {choice_1}
- {choice_2}
- ... (only included if choices is provided)

Answer should be json of the form:

```
{  
  "answer": "<your answer>",&br/>  "reason": "<why>"  
}
```

Answer formatting guidelines:

{type_specific_guideline}

If the question cannot be answered based on the context, set answer to “non-answerable”.

QA Type-specific Guidelines

yes_no: Answer “yes” or “no” only.
which: Answer with the exact choice text from the given choices.
description: Provide a concise, categorical description. Focus on object types/categories rather than specific item names. Do NOT include exits, doors, or room connections. Valid options: [object, container, supporter, food, key].
where: Answer with the location/room name.
what: Answer concisely with the specific information requested.
default: Answer concisely.

D The comparison of answerable QA accuracy and non-answerable QA accuracy

We compare the accuracy of answerable QA and non-answerable QA on the T2QBench dataset, as shown in Figure 8. In another word, we analyze the source of correctness. To figure out if agent can answer the exact fact ground truth, or gain accuracy simply by answering “I don’t know”.

From this figure, we observe that the composition of correctness sources varies substantially across models, particularly in the ratio of answerable vs. non-answerable QA. For location questions that can be resolved via single-turn retrieval, agents typically satisfy the prerequisite checkpoints during exploration and can answer correctly

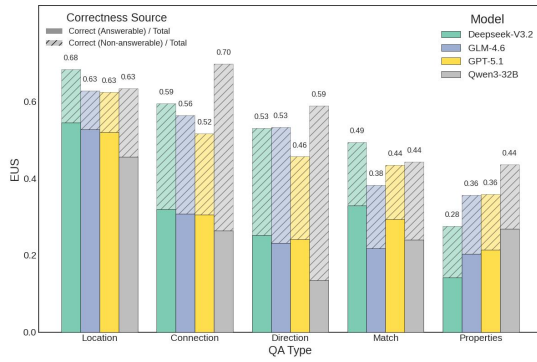


Figure 8: The comparison of answerable QA accuracy and non-answerable QA accuracy

913 once the relevant evidence is observed. In con-
 914 trast, for orientation and connectivity questions, in-
 915 sufficient coverage of key locations often makes
 916 the required evidence unavailable; consequently,
 917 models can frequently obtain “correct” outcomes
 918 by predicting non-answerable. For object-
 919 property questions, however, models perform
 920 poorly on both answerable and non-answerable in-
 921 stances, suggesting that limited proactive interac-
 922 tion (needed to reveal latent properties) remains a
 923 major bottleneck.