

Enter the Mind Palace: Reasoning and Planning for Long-term Active Embodied Question Answering

M. Fadhil Ginting^{1,2}, Dong-Ki Kim¹, Xiangyun Meng¹, Andrzej Reinke¹, Bandi Jai Krishna¹,
Navid Kayhani¹, Oriana Peltzer¹, David D. Fan¹, Amirreza Shaban¹, Sung-Kyun Kim¹,
Mykel J. Kochenderfer², Ali-akbar Agha-mohammadi¹, Shayegan Omidshafiei¹

¹Field AI, ²Stanford University
mind-palace-laeqa.github.io

Abstract: As robots become increasingly capable of operating over extended periods—spanning days, weeks, and even months—they are expected to accumulate knowledge of their environments and leverage this experience to assist humans more effectively. This paper studies the problem of Long-term Active Embodied Question Answering (LA-EQA), a new task in which a robot must both recall past experiences and actively explore its environment to answer complex, temporally-grounded questions. Unlike traditional EQA settings, which typically focus either on understanding the present environment alone or on recalling a single past observation, LA-EQA requires an agent to reason over past, present, and possible future states, deciding when to explore, when to consult its memory, and when to stop gathering observations and provide a final answer. Standard EQA approaches based on large models struggle in this setting due to limited context windows, absence of persistent memory, and an inability to combine memory recall with active exploration. To address this, we propose a structured memory system for robots, inspired by the mind palace method from cognitive science. Our method encodes episodic experiences as scene-graph-based world instances, forming a reasoning and planning algorithm that enables targeted memory retrieval and guided navigation. To balance exploration and recall, we introduce value-of-information-based stopping criteria that determine when the agent has gathered sufficient information. We evaluate our method on real-world experiments and introduce a new benchmark that spans popular simulation environments and actual industrial sites. Our approach significantly outperforms state-of-the-art baselines, yielding substantial gains in both answer accuracy and exploration efficiency.

Keywords: Embodied QA, long-term reasoning, vision-language navigation.

1 Introduction

Humans naturally develop long-term situational awareness through repeated interactions with their environment, remembering routines, recognizing object placements, and anticipating future needs. For example, when making a shopping list for breakfast, one can recall household preferences and check available supplies to identify what needs to be bought. This type of memory retrieval and long-term temporal grounding is key to intelligent embodied behavior. Among tasks related to this, Embodied Question Answering (EQA) is particularly compelling, as it probes a robot’s semantic understanding of its environment [1]. EQA approaches are typically framed either in active settings—where robots explore the environment from scratch to gather information [2]—or episodic settings—where robots answer questions using a single recorded trajectory [3]. While Vision-Language Models (VLMs) have improved performance [4, 5, 6], current approaches are limited to using only the robot’s present observations or a single episodic memory, and do not generalize to using multiple past experiences or long-term knowledge. To address this gap, we introduce

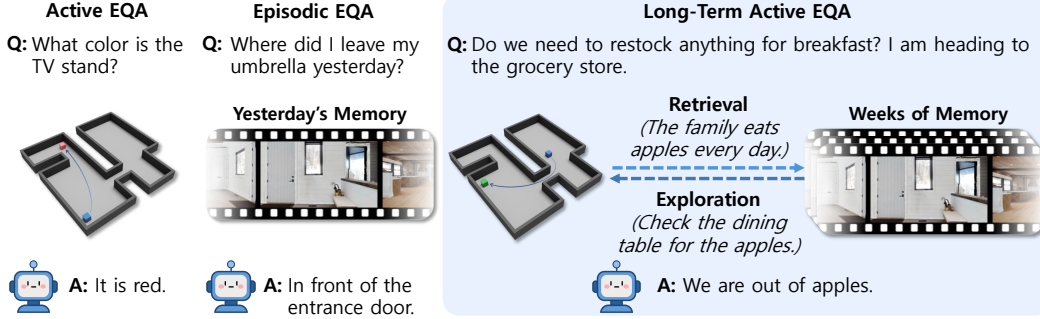


Figure 1: Different EQA problem setups. We study a new problem of Long-term Active EQA that combines active exploration with long-term memory understanding over multiple episodes.

Long-term Active Embodied Question Answering (LA-EQA), where robots must both recall past experiences and actively explore their surroundings to answer complex questions (see Figure 1). To our knowledge, this problem is largely unexplored, and no benchmark currently exists to evaluate it. Performing LA-EQA with VLMs and LLMs using existing EQA approaches is challenging for two reasons. First, representing the robot’s past observations accumulated over many deployments across days or months is difficult: a single run can generate thousands of images from diverse viewpoints, yet most questions only require a few relevant frames. Ingesting all this data directly is inefficient and often infeasible due to limited context windows. Second, retrieving relevant information from long-term memory and exploring relevant places in the environment creates a vast combined search space of past and new observations, where an uninformed search is computationally costly. These challenges raise a need for a new paradigm for long-term reasoning for embodied agents.

To address these challenges, we propose an approach for effective long-term memory representation and retrieval for embodied agents. Inspired by the mind palace technique [7]—where humans can effectively recall memories by associating them with spatial landmarks—we structure a robot’s long-term observations into a series of spatial world instances. Each instance is represented by a hierarchical scene graph that spatially groups semantic observations. Spatiotemporal structure is captured by linking multiple episodic world instances over time, enabling reasoning and exploration using retrieval of relevant experiences based on spatial proximity and temporal context. Our method, titled Mind Palace Exploration, has three components: 1) Generation, converting long-term memory into multiple scene-graph world instances; 2) Reasoning and Planning, where the robot interleaves EQA reasoning to identify target objects and assess if sufficient information has been gathered; and 3) Stopping Criteria, using *value of information* to balance memory recall and active exploration.

We introduce the first benchmark on LA-EQA and evaluate our approach against state-of-the-art baselines in EQA. In particular, the benchmark consists of diverse large-scale, high-fidelity simulation environments and real-world office and industrial sites across multiple days and months. Our approach outperforms baselines by 12–28% in answer correctness, achieves 16% higher exploration efficiency, and maintains a 12% correctness gain over the strongest baseline while using 77% fewer retrieved images, demonstrating both the effectiveness and efficiency of our approach. We further demonstrate the scalability and generalizability of our method in long-term settings across diverse environment types, including reasoning over memory from deployments spanning 2.4 km of robot trajectories collected over 6 months. We show the benefit of early memory retrieval stopping criteria in further reducing the number of past observation images while maintaining comparable performance. The real-world experiments demonstrate the feasibility of our approach in practical settings, where a legged robot deployed in a 1,000 m² office space uses past inspection memory to efficiently explore the environment and answer practical, day-to-day questions about the office.

2 Related Work

Embodied Question Answering (EQA) has been studied extensively from earlier works that employed learning-based models [1, 8, 9, 10, 11, 12] to more recent efforts leveraging foundation

models [3]. Recent approaches generally fall into two settings: episodic-memory EQA, where the agent accesses a single episode of memory, such as in OpenEQA [3] and ReMEmbR [13], and active EQA [14, 15, 16], where the agent explores a novel environment to gather information for answering questions, such as in Explore-EQA [2], Efficient-EQA [17], and Graph-EQA [18]. We propose a new and more general problem of Long-term Active EQA, in which the agent must integrate information across multiple prior episodes and active exploration to answer the question.

Semantic scene representation is a critical component for embodied reasoning and planning. Various methods have been proposed to encode the semantics and contextual structure of the world, including dense 3D representations [19, 20], voxel maps [21], and scene graphs [22, 23, 24]. In our work, we opt for a scene graph approach [25], which has demonstrated effectiveness in EQA tasks [26, 18, 27], and can be integrated with scalable memory retrieval and planning. We extend the scene graph from a single environment snapshot to a series of episodic scene graphs labeled by macro-temporal intervals (e.g., hours, days), enabling the agent to reason over multiple world instances that capture how the environment evolves across long-term deployments.

Semantic-guided navigation focuses on reasoning and planning methods for robot navigation directed by semantic cues, which has a rich body of literature [28, 29] involving tasks specified by images [30, 31], object categories [32, 33], and natural language [34, 35, 36, 37]. Our work related to semantic-based planning to search objects [38] and gather information for EQA tasks [17]. The problems are typically framed as either online planning, which builds representations incrementally during execution [39, 40, 41], or offline planning, which relies on pre-constructed maps of the environment [42, 43]. We address the challenge of leveraging multiple historical maps for online planning in long-term settings where the environment evolves over time. We propose a unified approach integrating offline memory retrieval with online exploration for LA-EQA.

3 Problem Formulation of Long-term Active EQA

LA-EQA is a setting where an agent answers questions about the environment by actively exploring it and retrieving relevant information from long-term memory. The LA-EQA task is defined as tuple (Q, M, E, x_0, A^*) , where Q is the question, $M = [m_1, \dots, m_N]$ is a list of episodic memories, E is the current environment, x_0 is the initial robot pose, and A^* is the ground truth answer. The environment is dynamic: its visual appearance and object states can change over time. Each episodic memory $m_i = [m_{i,1}, \dots, m_{i,L}]$ contains L tuples of past robot pose and image observations $m_{i,j} = (x_{i,j}, o_{i,j})$ collected within a specific macro-temporal interval (e.g., hours). This formulation is a generalization of active EQA and episodic memory EQA with a single episode (see Appendix A).

In LA-EQA, the agent follows policy $\pi(a_k | x_k, h_k, Q)$, mapping its state x_k at time step k , working memory h_k (history of action and observation since receiving Q), and the question to one of three possible actions: *retrieve*, *explore*, and *answer*. The *retrieve* action a^R recalls a past memory $m_{i,j}$ into h_k . The *explore* action a^E moves the robot to viewpoint w_i in E , storing the new observation o_k in h_k ; w_i need not be near the robot and can be any obstacle-free space informed by prior experience. The *answer* action a^T generates an answer A in natural language based on h and terminates the task.

4 Mind Palace Exploration for Solving Long-term Active EQA

Humans use the mind palace technique [7] to remember complex information by organizing it into a structured spatial memory, which enables efficient retrieval and traceable recall of relevant memories. We explore how this technique can be applied to long-term memory representation and reasoning in robots. Our approach consists of three key ideas (see Figure 2). First, prior to the EQA task, we construct a long-term memory representation (referred to as the Robotic Mind Palace \mathcal{M}), which summarizes the robot’s history of observations into multiple world instances of scene graphs $\mathcal{M} = [G_0, G_1, \dots, G_N]$. Then, during the LA-EQA scenario, the agent reasons over and explores these world instances in \mathcal{M} to answer the question Q using a policy π . Additionally, we introduce early stopping criteria using the notion of *value of information* to avoid retrieving memory that is unlikely to improve the next exploration action a^E . We provide a detailed algorithm in Appendix B.

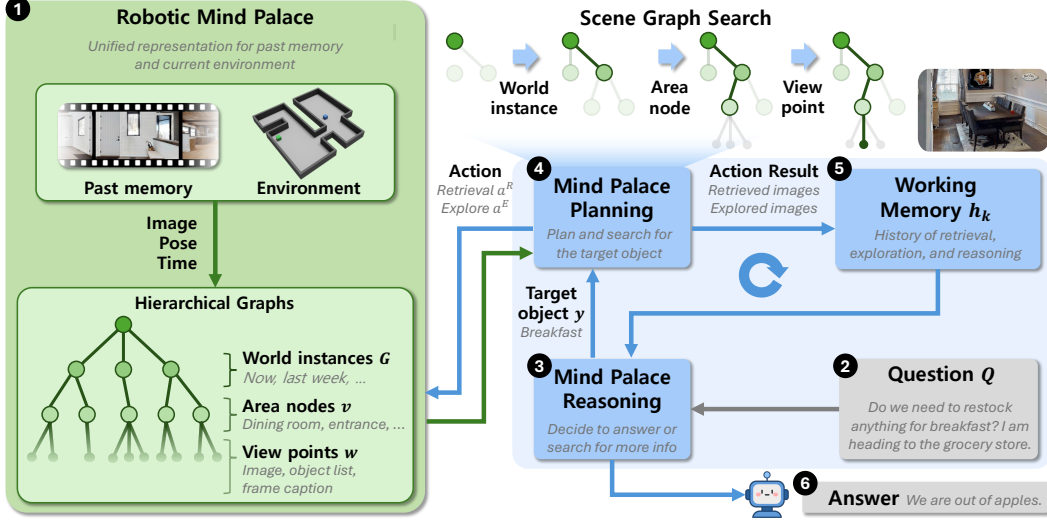


Figure 2: **Mind Palace Exploration** builds a Robotic Mind Palace that unifies past memories and environment representation (1). Given a question (2), the agent alternates between reasoning over the question to identify a target object (3), planning a search strategy through memory retrieval and exploration (4), and updating its working memory (5), until it is ready to answer the question (6).

4.1 Mind Palace Generation

Mind Palace is a series of episodic world instances. The Mind Palace divides the long-term history of robot image observation and trajectories M into episodes m based on a macro-temporal term such as hours, times of day, and weeks. The chunking of the episodes comes naturally in robotics as a mobile robot in continuous operations needs to pause any activities while recharging the battery. Each episode becomes a world instance in the Mind Palace and is indexed by its macro-temporal label in texts, allowing an LLM-based agent to select relevant episodes to recall.

An Episodic world instance is represented as a hierarchical scene graph. Given the sequence of robot observation and trajectory within an episode m_i , we build a world representation as a hierarchical scene graph $G_i = (\mathcal{V}_i, \mathcal{E}_i)$, where \mathcal{V}_i denotes the set of nodes and \mathcal{E}_i denotes the edges connecting the nodes [23]. First, we sample dense viewpoints w from the past trajectory to form a set of viewpoint nodes. Each viewpoint node w_i is associated with the robot pose x , images, a list of detected objects in the image [44, 45], and frame captions. The list of objects and frame captions is used as an index for LLM-based agents for image retrieval selection. Then the viewpoint nodes w are clustered into area nodes $v \in \mathcal{V}_i$ based on the spatial and contextual similarity [46, 26, 47]. Each area node v is associated with the centroid of all the clustered viewpoints and the object list. The neighboring viewpoints w and areas v are connected with graph edges, and every w is connected to a v , forming a hierarchical scene graph representation G for each world instance.

The Robotic Mind Palace consists of a series of world instances representing the *past* long-term memory $[G_1, \dots, G_N]$ and the *present* knowledge of the environment G_0 . At the start of the LA-EQA task, we assume the robot has not explored the present environment yet, so world instance G_0 is only initialized with area nodes v because the state of the environment and object placement may have changed since the last mapping in G_1 . We update G_0 as the robot explores the environment.

4.2 Mind Palace Reasoning and Planning

We perform reasoning and planning over the robotic mind palace to solve the LA-EQA task. This involves three interleaving steps: 1) reasoning over the question to determine what object or spatial concept y to search and when the agent can answer the question, 2) hierarchical planning over the Mind Palace to gather information, and 3) updating the information to the working memory h .

Reasoning over the LA-EQA: The first step in the reasoning process is to determine whether the robot has sufficient information to answer question Q using working memory h_k , which stores past actions, observations, and prior reasoning steps. The agent queries a VLM with h_k and Q . If the VLM responds it is possible to answer the question, the agent executes the *answer* action and provides an answer A with a VLM query. Otherwise, the agent queries an LLM to identify a target object or a spatial concept y , either a specific object explicitly stated in the question or an inferred cue (e.g., something to make a coffee), which becomes the next object goal for exploration.

Planning over episodic world instances G in the Mind Palace: Mind Palace planning begins by selecting a sequence of world instances G to locate y efficiently. We query an LLM with a two-step reasoning process because we observe direct query often yields inefficient plans. The first step asks the LLM to reason whether answering the question requires object search across multiple world instances or if it only concerns a specific instance. Based on this reasoning, the LLM selects a subset of $G \in \mathcal{M}$ and plans a sequence of G . We guide the sequential planning with a heuristic that suggests prioritizing past world instances over the present instance G_0 as using prior knowledge of y locations in the past can inform and improve object search efficiency in the present.

Planning over areas v in the scene graph: Given a world instance G_i , we plan a sequence of areas v to explore that maximizes the probability of finding y . This is framed as an object-goal navigation problem, and we adopt the planning formulation of object search over a scene graph [33]. We first query an LLM to output the probability of finding object y on each area $v \in G_i$, then use a forward search planner to find the best sequence of areas v to explore that minimizes the cost J to find y [48]. When exploring the present scene graph G_0 , the cost is defined by the path length between the robot’s current pose x_k and the centroid of each area. In contrast, when reasoning over past graphs, the agent can teleport to any area at a constant cost, regardless of the travel distance.

Exploring viewpoints w and replanning: Given an area v_i to search, we query the LLM to select viewpoints w based on the textual information in G_i . The object y may appear in frame captions or the object list but often is not mentioned, and relevant viewpoints must be inferred given the textual information [49]. The robot then explores the viewpoints by recalling images from the Mind Palace or navigating to the viewpoints in the environment using a robot-specific motion planner and taking the images. The retrieved or observed images are then stored in the working memory h_k . We repeat the planning over areas v and viewpoints w until the object y is detected in images by a VLM or until we reach the exploration limits. If the object is detected, we search for y in remaining world instances G and move to the *reasoning over the LA-EQA* step.

4.3 Early Stopping of Memory Retrieval for Navigation

This section examines how to reduce memory retrieval while maintaining exploration efficiency comparable to that of the unlimited memory retrieval case. In particular, we develop stopping criteria that decide when to halt past memory retrieval and proceed with exploration. Given a sequence of world instances that includes the present instance G_0 (e.g., $[G_1, G_2, G_0]$), we use an LLM to form a prediction set of areas $v \in G_0$, where the object y can be located with a probability above a threshold $P(y) \geq 1 - q$. Studies have shown that the LLM prediction and threshold $1 - q$ can be calibrated [2, 50, 51]. Using the prediction set, we define two possible conditions to immediately stop memory retrieval from past world instances $[G_1, G_2]$: 1) the prediction set contains only one area; 2) further memory retrieval will not improve the robot plan over the next sequence to explore v_i in the prediction set. We evaluate the possible improvements on the sequence using the notion of Value of Information (VoI) [52], which quantifies the expected utility gain from retrieving past memory, reducing the expected exploration cost J . We discuss more details in [Appendix B](#).

5 Long-term Active EQA Benchmark

Existing EQA datasets [3, 2, 13, 14] focus on scene understanding over short time spans (i.e., the same day), limiting their ability to capture long-term evolution of a scene (e.g., days and months). To address this, we curate the first LA-EQA dataset and benchmark, consisting of 3 simulated and 2



Figure 3: **LA-EQA Benchmark**: simulated and real-world scenes spanning multiple days / months.

real-world scenes (see Figure 3). For each simulation scene, we generate 5–10 scene variations over multiple days, reflecting changes caused by common routines. For real-world scenes, we collected 11 trajectories (30–60 mins) in an industrial site and an office environment over a 6-month period.

Question types: We categorize the questions based on their required temporal reasoning to capture different aspects of long-term scene understanding. **1) Past questions** pertain to a specific event observed in a single past trajectory. **2) Present questions** require only exploration of the current environment. **3) Multi-past questions** involve synthesizing information from multiple past trajectories (e.g., “What do we usually eat for breakfast?”). **4) Past-present questions** require reasoning over both historical memory and the current scene (e.g., “Are we missing anything we usually have for breakfast?”). **5) Past-present-future questions** involve predicting future outcomes based on both past and present observations (e.g., “When do you think we will run out of apples for breakfast?”).

We curated 150 questions, which uniformly cover the question types. The questions were generated by seven people to ensure the diversity of the questions. The dataset consists of past trajectories and observations, simulation environments, ground truth answers, and exploration solutions. More detailed statistics and representative questions are presented in Appendix C.

6 Experiments and Discussion

To evaluate our method, we answer: **Q1)** Does Mind Palace Exploration outperform other EQA methods across question types and lengths of past memory in long-term active EQA? **Q2)** Does early stopping of memory retrieval reduce the amount of memory retrieved without sacrificing performance? **Q3)** Can Mind Palace Exploration be practically deployed in real-world settings?

Methods: We compare our approach against the following baselines: **1) Multi-Frame VLMs** process the question with images and robot poses through a VLM to output the answer. This method is the strongest approach in the OpenEQA benchmark. **2) Socratic LLMs w/ Frame and Scene Graph Captions** use image and scene-graph captions and robot poses to answer the question. **3) Re-MEMBR [13]** is a state-of-the-art method in episodic EQA by building a queryable vector database representation of the robot pose, observation time, and image caption embedding and retrieving relevant entries in the database using an LLM. We use the open-source code of the method. **4) Active EQA Agent w/ Frames as the Memory** has the same information as Multi-Frame VLMs, but it lets the agent explore the environment by providing a list of viewpoints that the robot can visit. This approach is similar to the state-of-the-art method of using long-context VLMs with topological graphs [34] applied to the LA-EQA setting. **5) Active Socratic EQA Agent w/ Captions as the Memory** uses the same past memory information as Socratic LLMs w/ Frame and Scene Graph Captions, but it lets the agent explore viewpoints and analyze explored images with VLMs. All approaches use the GPT-4o as the language and vision model [5] and have the same maximum budget image retrieval, wherever applicable, and have the same exploration budgets on all active methods. The implementation details are provided in Appendix D.

Metrics: We evaluate all the agents using three metrics: **1) Answer correctness** is compared to the human-annotated answer and judged by an LLM-based scoring [3]. **2) Exploration efficiency** measures the path length of robot exploration compared to the oracle path length weighted by *answer correctness*. **3) Memory retrieval efficiency** measures the number of past images retrieved to answer the question. The details of the metrics are provided in Appendix E.

Methods	Answer	Expl. Eff.	Mem. (#)
Mind Palace (Ours)	65.0%	0.45	22.86
Mind Palace w/ stopping	61.8%	0.42	15.73
Multi-Frame VLMs [3]	52.9%	-	100
Socratic LLMs [3]	44.3%	-	0
ReMemBR [13]	46.1%	-	0
Active EQA w/ Frames	43.7%	0.29	100
Active Socratic EQA	36.8%	0.19	0

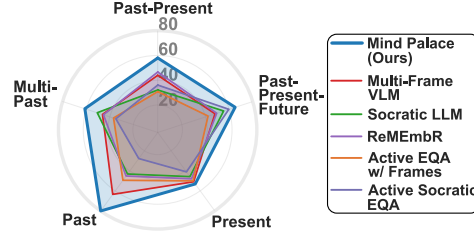


Table 1: LA-EQA results over answer correctness, exploration and retrieval efficiency. Figure 4: Performance over temporal reasoning question types.

6.1 Q1) Mind Palace Exploration Outperforms other EQA Approaches

Mind Palace Exploration outperforms baselines in all metrics. As shown in Table 1, our approach considerably outperforms all methods across the metrics, highlighting the gap in the current EQA approaches in the long-term EQA setting. More detailed results are provided in Appendix F.

Efficient past image retrieval is the key to multi-episodic world understanding. Our approach significantly outperforms the others that require specific information from past memory, represented by *past* and *multi-past* question types in Figure 4. This is largely because images convey richer visual contexts than captions, enabling more accurate answers about object properties, states, and placements. In the LA-EQA setting, multi-frame VLMs struggle as the maximum context length of the state-of-the-art VLMs is not comparable to the sheer amount of past observations in the memory. Our image retrieval approach is critical for efficient image analysis, as EQA questions typically need only several question-related images across multiple episodic memories. The results in Table 1 show that our approach only needs 77.14% fewer images compared to VLM-based methods, with much higher answer correctness.

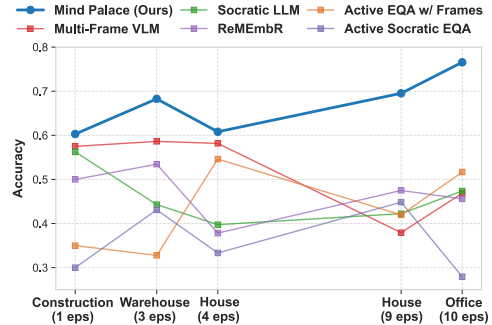


Figure 5: Performance of five different environments in the LA-EQA benchmark with a varying number of episodes.

Leveraging long-term memory improves active exploration efficiency. Our method achieves higher exploration efficiency than other active EQA agents (Table 1), particularly on *past-present* questions (Figure 4) in which past information can benefit present exploration. Our approach often recalls past memories to locate objects of interest so it can more accurately predict the probabilities of the object placements across areas v in the present environment.

Mind Palace Exploration is a scalable approach for LA-EQA. We evaluate scalability by plotting answer accuracy across different environments with varying numbers of past episodic memories (Figure 5). Our approach shows increasing performance gains over other methods with the number of past episodes in the memory. Given the same image retrieval limits, multi-frame VLM performance considerably drops as the images have less coverage across all the memories. ReMemBR performs steadily, highlighting the value of retrieval-based approaches in long-term EQA problems.

Our approach generalizes to diverse environments beyond the standard house setting. To test our approach further beyond standard EQA home environments benchmarks, we evaluated Mind Palace Exploration in larger real-world construction sites, a large office, and a simulated warehouse (Figure 5), where it consistently outperforms others, highlighting its flexibility. Building a structured memory representation for efficient exploration and retrieval becomes more critical as the environment size increases across many episodes. See Appendix I for full results.

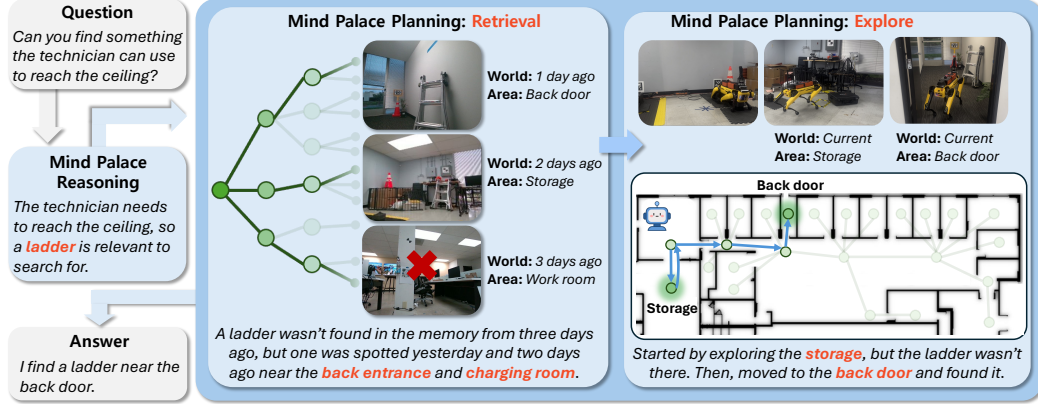


Figure 6: **Long-term Active EQA Hardware Experiments.** The robot retrieve relevant past information from the past memory and navigate around the office to answer the question.

6.2 Q2) Benefits of Early Memory Retrieval Stopping

Early memory retrieval stopping reduces the number of memories retrieved without sacrificing performance. As shown in Table 1, early stopping reduces the amount of image retrieval from the past memory while maintaining comparable answer accuracy. The early stopping reduces the number of past world instances that the agent retrieves if there is no new observation that will change the agent’s next exploration action. Examples in the experiment that we observe include when the robot predicts the possible areas where the object of interest is on the second floor, the robot will stop retrieving past world instances and move to the second floor. The stopping criteria are beneficial to even further improve the memory retrieval efficiency in Mind Palace Exploration.

6.3 Q3) Real-world Hardware Experiments

We demonstrate the efficacy of Mind Palace Exploration in real-world LA-EQA use cases in an office space spanning over 1,000 m² with 27 different areas, using a legged robot as an office assistant. The robot accesses 10 past episodes of past runs, inspecting the office for the past four days and six monthly inspections from October 2024 to March 2025. All the Mind Palace memory storage and planning, other than the GPT4-o query, is performed on the robot. A user sends the question to the robot remotely through a computer, and the robot reports back the answer once it finishes the task. We select 7 questions from the LA-EQA benchmark that require active exploration (Figure 6).

Mind Palace Exploration enables efficient exploration for practical real-world tasks. By consolidating knowledge of past object placements, the robot can efficiently locate relevant objects, saving an average of 3–10 room searches across the seven evaluated questions compared to a robot without memory access. The questions reflect realistic office scenarios (e.g., searching for tools, tracking missing packages, or identifying vacant desks unused for days) demonstrating the practical utility of LA-EQA. The robot can answer all the questions given that the information is available in its past memory and the current environment. See Appendix G for hardware experiment details.

7 Conclusion

We present the problem of LA-EQA, a new task that requires robots to combine long-term environment understanding with active exploration. We propose Mind Palace Exploration to address LA-EQA by representing past long-term memory and the present environment with a robotic mind palace, enabling reasoning and planning over the Mind Palace. We introduce the first benchmark for long-term active EQA, spanning days of simulation environments and months of real-world data, to foster future research in long-term reasoning. Our approach significantly outperforms state-of-the-art EQA baselines, highlighting the need for a new paradigm for LA-EQA.

8 Limitations

We discuss the limitations and challenges that we faced while evaluating our approach on the LA-EQA problem and building the Long-term EQA benchmark.

Long-term understanding is limited to the coverage of past trajectories. When evaluating our approach in simulations, we design questions that a human could answer by reviewing past memories or by exploring the environment. Our approach also assumes that questions can be answered using a combination of exploration and retrieval from past memories. However, in real-world hardware experiments, we find a question that requires broader past coverage to be answered accurately. The question asks the robot to find a vacant desk in the office that had not been used in recent days. Our approach mistakenly selects an occupied desk because, in the robot’s past data covering a one-hour walk every day, the desk is observed as vacant. This limitation highlights an exciting future research direction where robotic agents can recognize when they lack sufficient information and request additional information from humans or other agents.

Answer correctness gap due to limitations in semantic understanding of images and spatial understanding over the memories and environment. We observe that the cases where our approach does not achieve full scores on benchmark questions are primarily due to limitations of the VLMs’ semantic and spatial understanding capabilities. The semantic understanding issues include failures in detecting, counting, and understanding the functional properties of objects. Stronger VLMs or multiple specialized vision models can be easily integrated into the Mind Palace Exploration framework, and we expect performance to improve as the VLMs’ capabilities increase. Moreover, we observe several failed examples involving spatial understanding of images. For example, a VLM often assumes all objects detected while the robot is standing in one area are associated with the same area. This assumption does not hold when the images contain views of other areas. We find the most challenging questions for all approaches often require agents to collect multiple pieces of information and search different objects within the same environment and across multiple episodes. The problem highlights future research directions to further improve spatiotemporal reasoning capabilities for robots in a long-term setting, as measured by the performance in Long-term Active EQA benchmarks. We provide examples of the failure cases in [Appendix H](#).

Long-term EQA benchmarks require manual design by humans, limiting scalability. Our benchmark is currently designed by a team of human experts, and the team carefully curates the questions. However, this approach limits the scalability of the benchmark generation. We find several challenges in building the first long-term EQA benchmarks. First, in the simulation environment design, a human must carefully curate long-term scenarios spanning multiple days and determine the changes in object placements over time. We also need to place the objects manually, as we did not find reliable tools that can automate this process across different simulation worlds. Moreover, the questions still need to be carefully designed by humans to ensure their quality and relevance to typical long-term EQA tasks in home and industrial use cases. This challenge becomes more pronounced when creating questions from extensive real-world data. We experimented with using LLMs for question generation, but the results were still far from effective. Addressing these limitations in long-term EQA benchmark generation will require further research on automated scenario generation, scene modification, and practical long-term EQA question generation.

References

- [1] A. Das, S. Datta, G. Gkioxari, S. Lee, D. Parikh, and D. Batra. Embodied question answering. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [2] A. Z. Ren, J. Clark, A. Dixit, M. Itkina, A. Majumdar, and D. Sadigh. Explore until confident: Efficient exploration for embodied question answering. *Robotics: Science and Systems*, 2024.
- [3] A. Majumdar, A. Ajay, X. Zhang, P. Putta, S. Yenamandra, M. Henaff, S. Silwal, P. McVay, O. Maksymets, S. Arnaud, et al. OpenEQA: Embodied question answering in the era of foundation models. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- [4] H. Liu, C. Li, Y. Li, and Y. J. Lee. Improved baselines with visual instruction tuning. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- [5] A. Hurst, A. Lerer, A. P. Goucher, A. Perelman, A. Ramesh, A. Clark, A. Ostrow, A. Welihinda, A. Hayes, A. Radford, et al. GPT-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
- [6] A. Grattafiori, A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [7] E. L. Legge, C. R. Madan, E. T. Ng, and J. B. Caplan. Building a memory palace in minutes: Equivalent memory performance using virtual versus conventional environments with the method of loci. *Acta Psychologica*, 141(3):380–390, 2012.
- [8] D. Gordon, A. Kembhavi, M. Rastegari, J. Redmon, D. Fox, and A. Farhadi. IQA: Visual question answering in interactive environments. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [9] A. Das, G. Gkioxari, S. Lee, D. Parikh, and D. Batra. Neural modular control for embodied question answering. In *Conference on Robot Learning (CoRL)*, 2018.
- [10] E. Wijmans, S. Datta, O. Maksymets, A. Das, G. Gkioxari, S. Lee, I. Essa, D. Parikh, and D. Batra. Embodied question answering in photorealistic environments with point cloud perception. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [11] L. Yu, X. Chen, G. Gkioxari, M. Bansal, T. L. Berg, and D. Batra. Multi-target embodied question answering. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [12] J. Thomason, D. Gordon, and Y. Bisk. Shifting the baseline: Single modality performance on visual navigation & QA. In *North American Association for Computational Linguistics (NAACL)*, 2019.
- [13] A. Anwar, J. Welsh, J. Biswas, S. Pouya, and Y. Chang. Remembr: Building and reasoning over long-horizon spatio-temporal memory for robot navigation. *IEEE International Conference on Robotics and Automation (ICRA)*, 2025.
- [14] K. Jiang, Y. Liu, W. Chen, J. Luo, Z. Chen, L. Pan, G. Li, and L. Lin. Beyond the destination: A novel benchmark for exploration-aware embodied question answering. *arXiv preprint arXiv:2503.11117*, 2025.
- [15] V. S. Dorbala, P. Goyal, R. Piramuthu, M. Johnston, R. Ghanadhan, and D. Manocha. S-eqa: Tackling situational queries in embodied question answering. *arXiv preprint arXiv:2405.04732*, 2024.

- [16] T. Wu, C. Zhou, Y. H. Wong, L. Gu, and J. Yang. NoisyEQA: Benchmarking embodied question answering against noisy queries. *arXiv preprint arXiv:2412.10726*, 2024.
- [17] K. Cheng, Z. Li, X. Sun, B.-C. Min, A. S. Bedi, and A. Bera. Efficienteqa: An efficient approach for open vocabulary embodied question answering. *arXiv preprint arXiv:2410.20263*, 2024.
- [18] S. Saxena, B. Buchanan, C. Paxton, B. Chen, N. Vaskevicius, L. Palmieri, J. Francis, and O. Kroemer. Grapheqa: Using 3d semantic scene graphs for real-time embodied question answering. *arXiv preprint arXiv:2412.14480*, 2024.
- [19] S. Peng, K. Genova, C. Jiang, A. Tagliasacchi, M. Pollefeys, T. Funkhouser, et al. Openscene: 3d scene understanding with open vocabularies. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [20] N. M. M. Shafiullah, C. Paxton, L. Pinto, S. Chintala, and A. Szlam. Clip-fields: Weakly supervised semantic fields for robotic memory. In *ICRA Workshop on Pretraining for Robotics (PT4R)*, 2023.
- [21] P. Liu, Z. Guo, M. Warke, S. Chintala, C. Paxton, N. M. M. Shafiullah, and L. Pinto. Dynamem: Online dynamic spatio-semantic memory for open world mobile manipulation, 2024.
- [22] I. Armeni, Z.-Y. He, J. Gwak, A. R. Zamir, M. Fischer, J. Malik, and S. Savarese. 3d scene graph: A structure for unified semantics, 3d space, and camera. In *International Conference on Computer Vision (ICCV)*, 2019.
- [23] A. Rosinol, A. Violette, M. Abate, N. Hughes, Y. Chang, J. Shi, A. Gupta, and L. Carlone. Kimera: From SLAM to spatial perception with 3d dynamic scene graphs. *The International Journal of Robotics Research*, 40(12-14):1510–1546, 2021.
- [24] A. Werby, C. Huang, M. Büchner, A. Valada, and W. Burgard. Hierarchical open-vocabulary 3d scene graphs for language-grounded robot navigation. In *Robotics: Science and Systems*, 2024.
- [25] K. Rana, J. Haviland, S. Garg, J. Abou-Chakra, I. Reid, and N. Suenderhauf. Sayplan: Grounding large language models using 3d scene graphs for scalable robot task planning. *Conference on Robot Learning (CoRL)*, 2023.
- [26] Q. Xie, S. Y. Min, T. Zhang, K. Xu, A. Bajaj, R. Salakhutdinov, M. Johnson-Roberson, and Y. Bisk. Embodied-RAG: General non-parametric embodied memory for retrieval and generation. *arXiv preprint arXiv:2403.00000*, 2024.
- [27] Y. Yang, H. Yang, J. Zhou, P. Chen, H. Zhang, Y. Du, and C. Gan. Snapmem: Snapshot-based 3d scene memory for embodied exploration and reasoning. *arXiv preprint arXiv:2411.17735*, 2024.
- [28] P. Anderson, A. Chang, D. S. Chaplot, A. Dosovitskiy, S. Gupta, V. Koltun, J. Kosecka, J. Malik, R. Mottaghi, M. Savva, et al. On evaluation of embodied navigation agents. *arXiv preprint arXiv:1807.06757*, 2018.
- [29] M. Deitke, D. Batra, Y. Bisk, T. Campari, A. X. Chang, D. S. Chaplot, C. Chen, C. P. D’Arpino, K. Ehsani, A. Farhadi, et al. Retrospectives on the embodied ai workshop. *arXiv preprint arXiv:2210.06849*, 2022.
- [30] Y. Zhu, R. Mottaghi, E. Kolve, J. J. Lim, A. Gupta, L. Fei-Fei, and A. Farhadi. Target-driven visual navigation in indoor scenes using deep reinforcement learning. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2017.

- [31] L. Mezghan, S. Sukhbaatar, T. Lavril, O. Maksymets, D. Batra, P. Bojanowski, and K. Alahari. Memory-augmented reinforcement learning for image-goal navigation. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2022.
- [32] N. Yokoyama, S. Ha, D. Batra, J. Wang, and B. Bucher. Vlfm: Vision-language frontier maps for zero-shot semantic navigation. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2024.
- [33] M. F. Ginting, S.-K. Kim, D. D. Fan, M. Palieri, M. J. Kochenderfer, and A. akbar Agha-mohammadi. SEEK: Semantic reasoning for object goal navigation in real world inspection tasks. In *Proc. of Robotics: Science and Systems*, 2024.
- [34] H.-T. L. Chiang, Z. Xu, Z. Fu, M. G. Jacob, T. Zhang, T.-W. E. Lee, W. Yu, C. Schenck, D. Rendleman, D. Shah, et al. Mobility vla: Multimodal instruction navigation with long-context vlms and topological graphs. *Conference on Robot Learning (CoRL)*, 2024.
- [35] Y. Fu, D.-K. Kim, J. Kim, S. Sohn, L. Logeswaran, K. Bae, and H. Lee. Autoguide: Automated generation and selection of context-aware guidelines for large language model agents. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 37, pages 119919–119948, 2024.
- [36] M. F. Ginting, D.-K. Kim, S.-K. Kim, B. J. Krishna, M. J. Kochenderfer, S. Omidshafiei, and A.-a. Agha-mohammadi. Saycomply: Grounding field robotic tasks in operational compliance through retrieval-based language models. *IEEE International Conference on Robotics and Automation (ICRA)*, 2024.
- [37] Y. Long, W. Cai, H. Wang, G. Zhan, and H. Dong. Instructnav: Zero-shot system for generic instruction navigation in unexplored environment. *Conference on Robot Learning (CoRL)*, 2024.
- [38] M. Khanna, R. Ramrakhya, G. Chhablani, S. Yenamandra, T. Gervet, M. Chang, Z. Kira, D. S. Chaplot, D. Batra, and R. Mottaghi. Goat-bench: A benchmark for multi-modal lifelong navigation. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- [39] S. Kim, A. Bouman, G. Salhotra, D. Fan, K. Otsu, J. Burdick, and A. Agha-mohammadi. PLGRIM: Hierarchical value learning for large-scale exploration in unknown environments. In *International Conference on Automated Planning and Scheduling*, 2021.
- [40] Z. Zhao, W. S. Lee, and D. Hsu. Large language models as commonsense knowledge for large-scale task planning. *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- [41] J. Chen, G. Li, S. Kumar, B. Ghanem, and F. Yu. How to not train your dragon: Training-free embodied object goal navigation with semantic frontiers. In *Robotics: Science and Systems*, 2023.
- [42] B. Chen, F. Xia, B. Ichter, K. Rao, K. Gopalakrishnan, M. S. Ryoo, A. Stone, and D. Kappler. Open-vocabulary queryable scene representations for real world planning. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2023.
- [43] Q. Gu, A. Kuwajerwala, S. Morin, K. M. Jatavallabhula, B. Sen, A. Agarwal, C. Rivera, W. Paul, K. Ellis, R. Chellappa, et al. Conceptgraphs: Open-vocabulary 3d scene graphs for perception and planning. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2024.
- [44] M. Zhang, K. Qu, V. Patil, C. Cadena, and M. Hutter. Tag map: A text-based map for spatial reasoning and navigation with large language models. *Conference on Robot Learning (CoRL)*, 2024.

- [45] Y. Zhang, X. Huang, J. Ma, Z. Li, Z. Luo, Y. Xie, Y. Qin, T. Luo, Y. Li, S. Liu, et al. Recognize anything: A strong image tagging model. *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [46] Y. Chang, N. Hughes, A. Ray, and L. Carlone. Hydra-multi: Collaborative online construction of 3d scene graphs with multi-robot teams. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2023.
- [47] N. Kayhani, B. McCabe, and B. Sankaran. Semantic-aware quality assessment of building elements using graph neural networks. *Automation in Construction*, 155:105054, 2023.
- [48] M. J. Kochenderfer, T. A. Wheeler, and K. H. Wray. *Algorithms for Decision Making*. MIT Press, 2022.
- [49] D. Shah, M. R. Equi, B. Osiński, F. Xia, B. Ichter, and S. Levine. Navigation with large language models: Semantic guesswork as a heuristic for planning. In *Conference on Robot Learning (CoRL)*, 2023.
- [50] A. Z. Ren, A. Dixit, A. Bodrova, S. Singh, S. Tu, N. Brown, P. Xu, L. Takayama, F. Xia, J. Varley, et al. Robots that ask for help: Uncertainty alignment for large language model planners. *Conference on Robot Learning (CoRL)*, 2023.
- [51] K. Tian, E. Mitchell, A. Zhou, A. Sharma, R. Rafailov, H. Yao, C. Finn, and C. Manning. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2023.
- [52] R. A. Howard. Information value theory. *IEEE Transactions on Systems Science and Cybernetics*, 2(1):22–26, 2007.

Project Webpage

Our project website can be accessed at: <https://mind-palace-laeqa.github.io/>

Appendix Contents

Appendix A: Extended Problem Formulation of LA-EQA	15
Appendix B: Mind Palace Exploration Details	15
Appendix C: Long-term Active EQA Benchmark	21
Appendix D: EQA Methods Details	26
Appendix E: Evaluation Metrics	27
Appendix F: Example Results from Different Methods	30
Appendix G: Hardware Experiment Details	32
Appendix H: Examples of Failure Cases	32
Appendix I: Example of Full Results	34

A Extended Problem Formulation of LA-EQA

The Long-term Active EQA (LA-EQA) problem is a generalized formulation that integrates active EQA (A-EQA) and episodic memory EQA (EM-EQA) [3], and extends them to support reasoning over long-term, multi-episodic memories.

Active EQA (A-EQA) [2, 17] is defined by the tuple (Q, E, x_0, A^*) , where Q is the question, E is the current environment, x_0 is the initial robot pose, and A^* is the ground truth answer. The robot has no knowledge of E and must actively explore the environment until it determines it has enough information to produce an answer A .

Episodic Memory EQA (EM-EQA) [3] provides the robot with a single episodic memory m_i , which contains a sequence of pre-recorded robot poses and image observations $m_{i,j} = (x_{i,j}, o_{i,j})$. The task is defined as the tuple (Q, m_i, A^*) , and the agent must answer the question based solely on the provided memory m_i .

Long-term Active EQA (LA-EQA) unifies the A-EQA and EM-EQA formulations by giving the robot access to both multiple episodic memories and the ability to explore the current environment. This setting better reflects real-world scenarios, where humans answer questions by combining present exploration with the recall of relevant past experiences. The past episodic memories may range from a single trajectory m_i to a sequence of N trajectories $M = [m_1, \dots, m_N]$ accumulated over multiple days, weeks, or even months.

B Mind Palace Exploration Details

B.1 Implementation Details of Mind Palace Generation

Hierarchical scene graph generation: We sample the robot’s past trajectory to ensure maximal coverage while minimizing redundant images stored in the long-term representation. In the simulation environment, viewpoints are sampled every three simulation time steps. For real-world data, we sample every three seconds and remove viewpoints with identical poses when the robot is stationary. We use GPT-4o to generate captions for the viewpoint images and to list the objects visible in each image. The viewpoints are clustered into areas by querying VLMs using the viewpoint positions, captions, and images. For real-world office and industrial data, where floor plan maps are available, we cluster viewpoints based on the areas defined in the floor plans.

B.2 Implementation Details of Reasoning over the LA-EQA

The following prompt is used to determine whether the robot has sufficient information to answer the question.

Prompt 1: Reasoning over LA-EQA to determine whether the robotic agent has sufficient information to answer the question

You are an AI agent operating in an environment. Your task is to answer user questions by either exploring the environment or recalling relevant past information.

The user has asked the following question:

Question: {question}

Here is a summary of the information collected so far from exploration or memory recall:

{working_memory}

Prompt 1 (continued)

Based on this information, do you think you have enough information to answer the question?
If not, explain why and identify which objects or areas still need to be explored.

Respond with a brief explanation to support your answer.

If the agent believes it still needs to gather more information, the following prompt template is used to identify the target object or spatial concept.

Prompt 2: Reasoning over LA-EQA to identify object of interest or spatial concept to search

To answer the user's question, what specific objects or places in the house should the agent focus on?
Specify an object, contextual description, or entity that could be observed by a robot's camera while walking around a house.

Question: `{user_question}`

For context, here is a summary of the agent's exploration and observations so far:
`{working_memory}`

It is important to explicitly include any object mentioned in the question within the answer. If the question does not mention a specific object but instead refers to an object by function or property (e.g., something that can be used for X or something with a certain color), respond with a more general object description rather than naming a specific item.

For instance, if the question is "What can I use to make tea?", respond with:
"something that can be used to make tea" | not with a specific answer like "kettle," since the available objects are still unknown.

Do not hypothesize the object's location at this stage. The goal is to identify objects or entities relevant to the question, not to guess where they are.

Keep the object description under 10 words.

If the agent believes it has sufficient information, the following query is used to return the answer.

Prompt 3: Reasoning over LA-EQA to answer the question

The user has asked the following question:
Question: `{question}`

We are now ready to answer the question based on the following working memory, which includes relevant observations and past exploration:
`{working_memory}`

Prompt 3 (Continued)

Based on this information, please answer the user's question directly. Provide a concise answer and reason to your answer.

B.3 Implementation Details Planning over episodic world instances

Benefits of the two-step reasoning process: The following example highlights the need of a two-step reasoning process to query LLM to plan over an episodic world instances.

Example 1: Benefits of the two-step reasoning process in planning over episodic world instances

Question: Do you know where is my backpack? I forget where I put it.

List of episodic world instances: ['now', 'yesterday evening']

Direct query: We need to search ['now'].

Two-step reasoning process: Recalling the placement of the backpack from the memory will help to search the backpack in the present. We need to search ['yesterday evening', 'now'].

We use the following prompt template for planning over episodic world instances G .

Prompt 4: Planning over episodic world instances

In this query, we aim to identify time instances|either through present exploration or past memory recall| that are most relevant for finding the target object $\{y\}$ in order to answer the user's question.

Currently, we are searching for the object: $\{y_object_to_search\}$

User question: $\{Question\}$

Reason for searching this object: $\{y_reasoning_to_search_object\}$

Available time instances: $\{list_of_world_instances_time\}$

We want to reason about which of the following search strategies is most efficient:

1. PAST_ONLY, Answerable using only past memory.
2. PRESENT_ONLY, Only present exploration is relevant.
3. PAST_THEN_PRESENT | Use past memory first to guide efficient current search.
4. MULTI_PAST_AND_PRESENT, Compare object state across time (e.g., trends, changes).

Ask yourself: If I want to find object $\{y\}$, does recalling past memory first help?

Prompt 4 (Continued)

Will it reduce time and energy needed to explore the present?
If the question involves the current state, consider both past and present.

Preferred strategy order: PAST_ONLY > PAST_THEN_PRESENT > PRESENT_ONLY > MULTI_PAST_AND_PRESENT
If selecting PAST_THEN_PRESENT or MULTI_PAST_AND_PRESENT, choose at most 5 past time instances.

Based on your selected strategy, propose the ordered list of time instances to search that are most relevant.

For context, here is the summary of the agent’s working memory and observations so far:

`{working_memory}`

B.4 Implementation Details Planning over Areas

Extended formulation of planning over areas v : We adopt the object search planning formulation over a scene graph [33] to minimize the expected total path length required to find the target object y . This problem is modeled as a Markov Decision Process (MDP), defined by the tuple (S, A, T, C) , representing the state space, action space, transition probabilities, and cost function, respectively.

The state space S consists of all possible areas $v \in G_i$ where the robot can be located, along with a terminal state s^y representing successful detection of object y . The action space A includes all areas the robot can explore. The transition function T includes both the transition probabilities between areas and the probabilities of transitioning to the goal state. The transition probabilities to the goal state are estimated using an LLM, which is queried with the following prompt.

The cost function C assigns a cost to each action, which corresponds to the path length from the robot’s current pose to the centroid of the target area, or the distance between areas. In the case of past image retrieval, each action incurs a constant cost.

We search for a sequence of areas $[v_i, v_{i+1}, \dots]$ that minimizes the expected cost-to-go J . In our implementation, we perform a forward search up to three steps ahead to select the next area to explore, and return the first area v_i in the optimal sequence.

We use the following prompt template to get the transition probabilities.

Prompt 5: LLM query to get transition probabilities to find the object y

You are an AI agent in an environment. Your task is to answer questions from the user by either exploring the environment or recalling past relevant information.

To locate the object: `{y-object-to-search}`, and to answer the question: `{Question}`, you must assess the probability (from 0.0 to 0.99) of finding the object in each area.

Assign a probability score to each area. Higher values mean greater confidence that the object is located in that area.

Only provide up to 10 areas with the highest probabilities. Do not include areas with very low likelihood.

Prompt 5 (Continued)

Here is the list of areas in the environment:
{list_of_area_names}

We are currently exploring the environment at world instance:
{G_i_world_instance_to_explore}

For context, here is the agent’s working memory across all world instances:
{working_memory}

B.5 Implementation Details Planning over Viewpoints

We use the following prompt template to get a list of viewpoints w to explore.

Prompt 6: Planning over viewpoints

To locate the object: {y_object_to_search}, and to answer the question: {Question}, what viewpoints should be searched? List at most five.

Here is information about the viewpoints in the selected area v_i.
Note: the objects listed for each place are not exhaustive|they only include easily identifiable items.
{list_of_viewpoints_in_area}

The list of the viewpoints are then passed to a classical path planner and in real-world experiment, passed to a legged-robot navigation policy.

B.6 Early Stopping of Memory Retrieval for Navigation

Preliminaries on Value of Information: The Value of Information (VoI) quantifies how much observing additional variables is expected to improve utility, or equivalently, reduce the expected cost-to-go J [52]. Let $J(o)$ denote the expected cost of an optimal action given observation o . Then, the VoI of observing a new variable O' is defined as:

$$VOI(O' | o) = J(o) - \sum_{o'} P(o' | o) J^*(o, o'). \quad (1)$$

This formulation captures the expected decrease in cost resulting from observing O' . If observing O' does not affect the optimal action, then its VoI is zero.

We apply this principle to define early stopping conditions for searching for the object y in the present environment G_0 , beginning with memory retrieval from past world instances (e.g., $[G_1, G_2, G_0]$).

Stopping Condition 1: The prediction set contains only one area: If the prediction set (formed using the transition probability estimate T , see Appendix B.4) contains only one area $v_i \in G_0$ where object y is likely to be found, and we assume y is present in the environment and the prediction set contains the true location of y , then no additional observation can change the target area. In this case, the VoI of recalling past memory is zero, as it does not influence the search strategy for y .

Stopping Condition 2: Further memory retrieval will not improve the next area to explore: Consider a prediction set v_1, v_2, v_3 , where v_1 is the closest area to the robot, and reaching v_2 or v_3 would require passing through v_1 . In this case, recalling additional past information will not change the robot’s next area to explore. Thus, the VoI of retrieving more memory is zero, as it does not help prioritize the next search area for y .

B.7 Mind Palace Exploration algorithm

Algorithm 1 Mind Palace Exploration

Require: Question Q , episodic memories M , environment E , initial robot pose x_0

```

1: # Step 1: Robotic Mind Palace generation
2:  $\mathcal{M} \leftarrow \text{mind\_palace\_generation}(M, E)$ 
3: # Step 2: Working Memory and state initialization
4:  $h_k \leftarrow []$ 
5:  $x_k \leftarrow x_0$ 
6:  $a_k \leftarrow \text{null}$ 
7: # Step 3: Loop until ready to answer
8: repeat
9:   # Mind Palace Reasoning
10:  if  $\text{is\_ready\_to\_answer}(h_k, Q)$ :
11:     $a_k \leftarrow a^T$ 
12:    break
13:   $y \leftarrow \text{identify\_target\_object}(h_k, Q)$ 
14:  # Mind Palace Planning
15:   $G \leftarrow \text{planning\_over\_world\_instances}(y, \mathcal{M}, Q, h_k)$ 
16:  for each  $G_i \in G$ :
17:     $v \leftarrow \text{planning\_over\_areas}(y, G_i, Q, h_k)$ 
18:    for each  $v_i \in v$ :
19:       $\text{bool\_object\_found} \leftarrow \text{False}$ 
20:       $w \leftarrow \text{planning\_over\_viewpoints}(y, v_i, Q, h_k)$ 
21:      for each  $w_i \in w$ :
22:        if  $G_i == G_0$ :
23:           $o_k \leftarrow \text{navigate}(w_i, E)$ 
24:           $\mathcal{M} \leftarrow \text{update\_mind\_palace}(o_k, \mathcal{M})$ 
25:        else:
26:           $o_k \leftarrow \text{retrieve}(w_i, \mathcal{M})$ 
27:           $\text{bool\_object\_found} \leftarrow \text{is\_y\_detected}(o_k)$ 
28:          if  $\text{bool\_object\_found}$ :
29:             $h_k \leftarrow \text{update\_working\_memory}(o_k, h_k, G_i, v_i, w_i)$ 
30:            break
31:      if  $\text{bool\_object\_found}$ :
32:        break
33:    if  $a_k == a^T$ :
34:      break
35: until  $a_k == a^T$ 
36: return  $\text{answer\_the\_question}(Q, h_k)$ 

```

C Long-term Active EQA Benchmark

The long-term active EQA benchmark consists of 150 questions annotated by a human team across five different environments. Each environment contains 30 questions that uniformly cover all five question types. This section present samples of the environment and the questions.

C.1 Habitat HM3D House 1



Figure 7: Sample images from the Habitat HM3D House 1.

Question	Type	Answer
What drink do we have at the kitchen counter now?	Present	Orange juice
Can you check if the package that arrived on Thursday afternoon is still located at the same place?	Past-present	Yes, the package is still on the same place
What happens if I leave the drink on the kitchen island for a few days? Can I still drink it after 3 days?	Past-present-future	The orange juice will be spoiled
Did I leave the orange juice on the kitchen counter overnight?	Past	No you did not
How long has the package been near the front door?	Multi-past	2 days since Thursday afternoon

Table 2: Example questions for the Habitat HM3D House 1.

#	Episodic Time Instance
1	Friday Afternoon
2	Thursday Afternoon
3	Thursday Morning
4	Wednesday Afternoon
5	Wednesday Morning
6	Tuesday Afternoon
7	Tuesday Morning
8	Monday Afternoon
9	Monday Morning
10	Sunday Afternoon

Table 3: List of past episodic memories from the Habitat HM3D House 1.

C.2 Habitat HM3D House 2

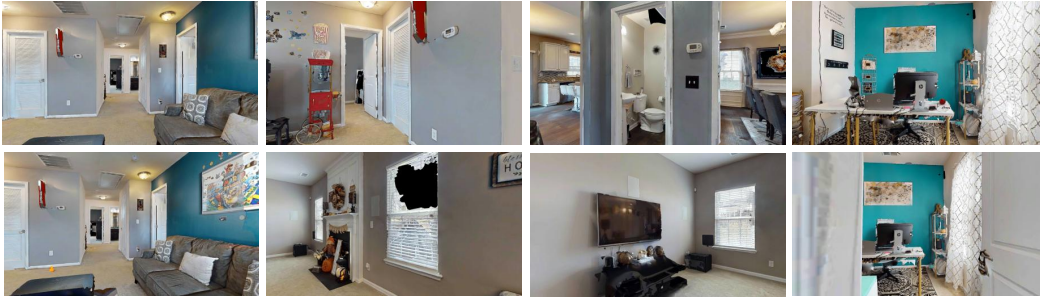


Figure 8: Sample images from the Habitat HM3D House 2.

Question	Type	Answer
What toys currently tucked on the stairways?	Present	A duck toy
I'm looking for a water bottle, where is it?	Past-present	It's still on the desk in the study room
I'm going to the grocery, what kind of fruit we should buy for the house	Past-present-future	Apples, I saw apples in different places
Where was the apple placed after the meal yesterday	Past	On the table in the Dining Room
Where have the kids left the duck toys for the past days?	Multi-past	On the floor in the living room and on the stair

Table 4: Example questions for the Habitat HM3D House 2.

#	Episodic Time Instance
1	Thursday Afternoon
2	Wednesday Afternoon
3	Tuesday Afternoon
4	Monday Afternoon

Table 5: List of past episodic memories from the Habitat HM3D House 2.

C.3 NVIDIA Isaac Large Warehouse

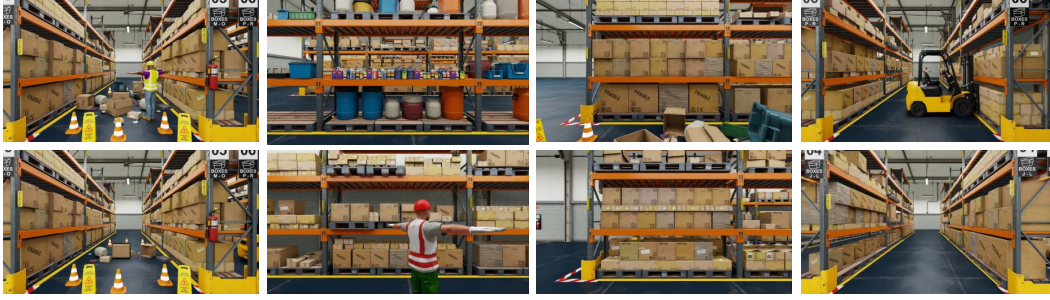


Figure 9: Sample images from the NVIDIA Isaac Large Warehouse.

Question	Type	Answer
How many blue barrels do we have on the shelves in Aisle 2?	Present	7
Can you check if the clutter has been partially cleaned up since yesterday?	Past-present	Yes, some clutter has been removed today
We want to recruit more workers in the warehouse. Which day should we assign them?	Past-present-future	There were no workers on Tuesday
Was there any person close to any hazards yesterday?	Past	Yes, in Aisle 5
Can you check if the fire extinguisher in Aisle 1 for the past days?	Multi-past	It's been placed at the same place

Table 6: Example questions for the NVIDIA Isaac Large Warehouse.

#	Episodic Time Instance
1	Wednesday
2	Tuesday
3	Monday

Table 7: List of past episodic memories from the NVIDIA Isaac large warehouse.

C.4 Real-world Office Environment

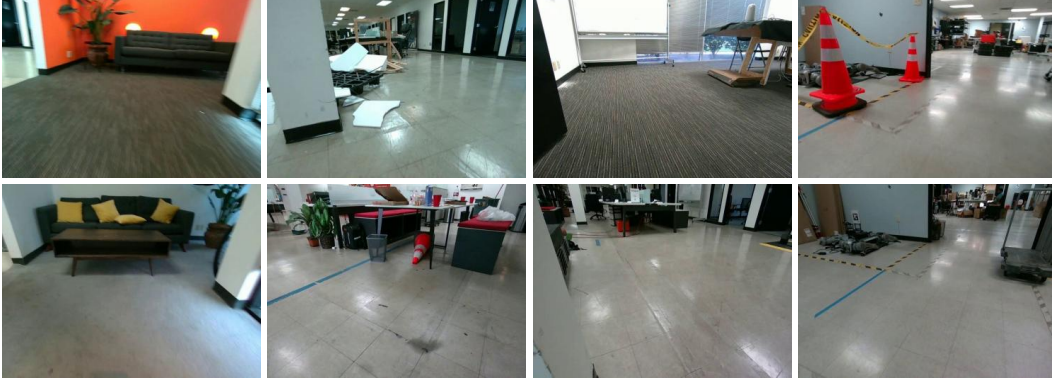


Figure 10: Sample images from the real-world office environment.

Question	Type	Answer
Is the door from the PPE room to outside closed or open?	Present	It's closed
The reception area wall of the office looks different from the one I last remembered in November, what's the difference?	Past-present	The wall is now painted orange
Can you check if the office side door through the PPE room is open? Do we usually keep it closed?	Past-present-future	We usually keep it open
Do you remember the type of the orange car parked on the side of the road in January?	Past	Land Rover
What is the fire extinguisher color that we placed in the robot room in December and January?	Multi-past	Red

Table 8: Example questions for the real-world office environment.

#	Episodic Time Instance
1	October 18, 2024
2	November 15, 2024
3	December 13, 2024
4	January 17, 2025
5	February 14, 2025
6	March 14, 2025
7	March 17, 2025 (Four days ago)
8	March 18, 2025 (Three days ago)
9	March 19, 2025 (Two days ago)
10	March 20, 2025 (Yesterday)

Table 9: List of past episodic memories from the real-world office environment.

C.5 Real-world High-rise Construction Site



Figure 11: Sample images from the real-world high-rise construction site.

Question	Type	Answer
What tools are on the wooden boxes?	Present	Drilling machines
Do you see the blue tape on the floor? Is it broken anywhere?	Past-present	Yes
What is the function of the blue tape used on the floor?	Past-present-future	To mark boundaries or safety zones
What was the color of the fan on the ground?	Past	White
Is the yellow cart loaded or empty?	Multi-past	Empty

Table 10: Example questions for the real-world high-rise construction site.

#	Episodic Time Instance
1	August 16, 2024

Table 11: List of past episodic memories from the real-world high-rise construction site.

D EQA Methods Details

Methods	Poses	Past Img.	Pres. Img.	Past Cap.	Pres. Cap.	SG Cap.	Expl.	Full Cov.
Mind Palace (Ours)	✓	✓	✓	✓	✗	✓	✓	✗
Multi-Frame VLMs [3]	✓	✓	✓	✗	✗	✗	✗	✓
Socratic LLMs [3]	✓	✗	✗	✓	✓	✓	✗	✓
ReMEMbR [13]	✓	✗	✗	✓	✓	✗	✗	✓
Active EQA w/ Frames	✓	✓	✓	✗	✗	✗	✓	✗
Active Socratic EQA	✓	✗	✓	✓	✗	✓	✓	✗

Table 12: **Comparison of information sources used by each method.** Poses = robot poses; Past Img. = access to past images for VLM analysis; Pres. Img. = access to present images for VLM analysis; Past Cap. = past image captions; Pres. Cap. = present image captions; SG Cap. = scene graph captions; Expl = supports active exploration; Full Cov. = requires a pre-recorded full-trajectory coverage before answering questions.

D.1 Multi-Frame VLMs

Prompt 7: Multi-Frame VLMs

You are an AI agent in an environment. Your task is to answer questions from the user by analyzing past image observations collected by the robot.

Question: {Question}

Use the following information to guide your reasoning and answer the question. Do your best to answer the question based on the information you have.

{images, robot poses}

Think step by step.

D.2 Socratic LLMs

Prompt 8: Socratic LLMs

You are an AI agent in an environment. Your task is to answer questions from the user by analyzing past image observations collected by the robot.

Question: {Question}

Use the following information to guide your reasoning and answer the question. Do your best to answer the question based on the information you have.

{image captions, scene graph captions, robot poses}

Think step by step.

D.3 ReMEMbR

We use the open-source code of ReMEMbR [13] that uses GPT4o for image captioning and LLM query. Following their real-world setup, ReMEMbR explores all the viewpoints in the present environment to build the memory prior to the question answering.

D.4 Active EQA w/ Frames

Prompt 9: Active EQA w/ Frames

You are an AI agent in an environment. Your task is to answer questions from the user by analyzing past image observations collected by the robot.

Question: {Question}

Use the following information to guide your reasoning and answer the question or if you think the question requires exploration of the present state of the environment, list some viewpoints for the robot to explore in the present environment to be able to answer the question.

{images, robot poses}

Think step by step.

D.5 Active Socratic EQA

Prompt 10: Active Socratic EQA

You are an AI agent in an environment. Your task is to answer questions from the user by analyzing past image observations collected by the robot.

Question: {Question}

Use the following information to guide your reasoning and answer the question or if you think the question requires exploration of the present state of the environment, list some viewpoints for the robot to explore in the present environment to be able to answer the question.

{image captions, scene graph captions, robot poses}

Think step by step.

D.6 Implementation Details and Parameters

All approaches and ours use the GPT-4o as the language and vision model [5] and for image captioning. In the experiments, the maximum budget of past image retrieval is 100 images and the maximum limit of the exploration is 25 viewpoints.

E Evaluation Metrics

We describe the three metrics used to evaluate robotic agents' performance in LA-EQA: 1) Answer correctness (Appendix E.1), 2) Exploration efficiency (Appendix E.2), and 3) Memory retrieval efficiency (Appendix E.3).

E.1 Answer correctness

The *answer correctness* metric evaluates how closely the agent's answer A matches the human-annotated answer A^* . Since answers to EQA questions are often open-ended, there may be no single

exact string that defines the correct answer, and multiple valid variations can exist. To account for this, we adopt the *LLM-Match* evaluation procedure introduced in the OpenEQA benchmark [3], which uses an LLM to score answer correctness.

This approach is faster and more cost-effective than manual human reviews and has demonstrated a high level of agreement with human judgments [3]. Given a question Q , a human-annotated answer A^* , and the agent’s answer A , the GPT-4o model is prompted to assign an integer score σ from 1 to 5, where 1 represents an incorrect answer and 5 represents a fully correct response.

The answer correctness score is then normalized to a percentage:

$$\mathcal{C} = \frac{\sigma - 1}{4} \times 100\%. \quad (2)$$

The overall answer correctness in the benchmark is the average of \mathcal{C} over all the questions.

Prompt 11: Prompt used for LLM-Match scoring

You are an AI agent. Your task is to evaluate the response given a question and the correct answer annotated by a human.

To mark the response, output a single integer from 1 to 5 (inclusive).

Question: {Question}

Correct answer: {GT_A_answer}

Response: {A_answer}

E.2 Exploration efficiency

The *exploration efficiency* metric measures how efficient the agent’s path to gather necessary information in the environment to answer the question correctly. We compare the total path length of the agent’s from the start of the LA-EQA task until it answers the question against a human-annotated solution. The agent’s total path length p is measured by calculating the total path distance from the starting pose x_0 through the sequence of the viewpoints w explored by the agent until it answers the question. The human annotated path length l is measured by calculating the total distance to a sequence of viewpoints w^* from x_0 that are annotated by humans that are sufficient to answer the questions.

Not all questions in the LA-EQA benchmark require present-environment exploration—for instance, questions that involve only past states. For such cases, the human-annotated sequence is set to empty, and $l = 0$. Furthermore, the human-annotated sequence w^* is not necessarily the only valid path to answer a question. For example, when annotating a solution for a question that requires checking objects on a dining table, a human may select the nearest viewpoint that provides a clear view, even though alternative, more distant or partially occluded viewpoints might also contain sufficient information. Thus, it is possible for the agent to find a shorter valid path than the annotated reference.

Given the agent’s total path length p and the human-annotated path length l , we adapt the Success weighted by (normalized) Path Length (SPL) metric—commonly used to evaluate navigation efficiency [28, 3]—to the LA-EQA setting:

$$\mathcal{X} = \begin{cases} \frac{\sigma - 1}{4} \times 100\%, & \text{if } l = 0 \text{ and } p = 0, \\ \frac{\sigma - 1}{4} \times \frac{l}{\max(l, p)} \times 100\%, & \text{otherwise.} \end{cases} \quad (3)$$

E.3 Memory Retrieval efficiency

The *memory retrieval efficiency* metric measures the number of images retrieved by the agent to answer a given question. Images generally contain richer information than captions for answering different type of questions. Measuring image retrieval efficiency is particularly important in long-term memory settings, as images represent larger amount of data that must be retrieved from memory storage and processed by a VLM compared to texts. We count the number of images retrieved after the agent receives the question.

F Example Results from Different Methods

F.1 Example results from different methods

Question: Do you know where my red backpack is? I forget where I left it. Correct answer: On the sofa in the living room downstairs.			
	Ours	Multi-frame VLM	Socratic LLMs
Retrieval and exploration	Retrieve 15 images and active exploration for 6.6 m	Retrieve 100 images and explore all viewpoints for 24.5 m	Explore all viewpoints for 24.5 m
Answer	Your red backpack is on the couch in the living room downstairs.	Your red backpack is likely in the entryway.	Your red backpack was last seen in the hallway on a bench.
Answer correctness	100%	25%	25%
Exploration Efficiency	54.3%	3.6%	3.6%
Images Retrieved	15 images	100 images	0 image
Reasoning	The exploration of the living room on Saturday morning shows the red backpack on the couch. This matches the past observation from Friday afternoon, confirming its current location.	The red backpack was last seen hanging in the entryway on Thursday afternoon. It was not visible in the images from Friday afternoon or Saturday morning, suggesting it might have been moved or is not in the current view. However, since it was last seen there, it is a good place to start looking.	The red backpack was observed multiple times in the hallway area, specifically on a bench, during the episodes from Thursday and Wednesday. It was not seen in the Saturday morning episode, suggesting it might still be in the same location where it was last observed.
	ReMEmbR	Active EQA w/ Frames	Active Socratic EQA
Retrieval and exploration	Explore all viewpoints for 24.5 m	Retrieve 100 images and explore viewpoints for 21.24 m	
Answer	Your red backpack is hanging on a hook in the hallway with wooden stairs and a black door.	I couldn't find your red backpack in the explored areas. It might have been moved to a different location.	Your red backpack was last seen in the hallway areas with benches and hooks.
Reasoning	Your red backpack is hanging on a hook in the hallway with wooden stairs and a black door.	The red backpack was last seen in the hallway near the stairs on Thursday afternoon. The images from Saturday morning show the hallway and the stairs, but the backpack is not visible in any of these images. This suggests that the backpack has been moved from its last known location. Since the images explored do not show the backpack, it is likely in a different area not covered by the current exploration.	The red backpack was seen in multiple places during the week, but most recently on Thursday afternoon, it was observed in several locations such as Place Node 18, Place Node 21, Place Node 24, Place Node 27, Place Node 30, Place Node 33, Place Node 36, and Place Node 39. Since the user mentioned they forgot where they left it last night, it is likely still in one of these locations.
Answer correctness	0%	25%	25%
Exploration efficiency	0%	4.2%	25%
Images Retrieved	0 image	100 images	0 image

Table 13: Example results from different methods

F.2 Example comparison on early stopping

Question: Do you know where my red backpack is? I forget where I left it. Correct answer: On the sofa in the living room downstairs.		
	Ours	Ours w/ early stopping
Retrieval and exploration	Retrieve 15 images (Friday afternoon, thursday afternoon, wednesday afternoon), and active exploration for 6.6 m	Retrieve 5 images (Friday afternoon), early stopping of image retrieval (not searching thursday afternoon, wednesday afternoon) , and active exploration for 5.7 m
Answer	Your red backpack is on the couch in the living room downstairs.	The red backpack is on the couch in the living room downstairs.
Reasoning	The exploration of the living room on Saturday morning shows the red backpack on the couch. This matches the past observation from Friday afternoon, confirming its current location.	The exploration of both past and present environments consistently shows the red backpack on the couch in the living room downstairs. The images from both Friday afternoon and Saturday morning confirm this location.
Answer Correctness	100%	100%
Exploration Efficiency	54.3%	61.8%
Images Retrieved	15 images	5 images

Table 14: Comparison of results with and without early stopping of memory retrieval

G Hardware Experiment Details

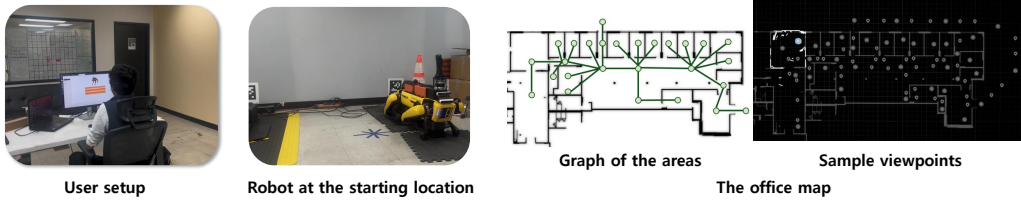


Figure 12: Hardware experiment setup in the office environment.

G.1 List of past episodic memories

#	Date	Distance (meters)
1	October 18, 2024	224
2	November 15, 2024	258
3	December 13, 2024	191
4	January 17, 2025	503
5	February 14, 2025	492
6	March 14, 2025	480
7	March 17, 2025	128
8	March 18, 2025	134
9	March 19, 2025	120
10	March 20, 2025	115

Table 15: List of past episodic memories collected from robot trajectories in the office.

G.2 List of questions

#	Questions used for the hardware experiments
1	Please inspect the fire extinguisher in the robot room and let me know if it has been in its usual location over the past few months.
2	Do you know where my jacket is? I think I left it in the meeting room yesterday, but I didn't see it there this morning.
3	Can you find a desk for the guest to use today? Look for one that hasn't been used in the past few days.
4	Can you check if the office side door through the PPE room is open? Do we usually keep it closed?
5	I'm looking for a package that was delivered. It's a large orange case. Do you know where it is?
6	We are changing the ceiling light bulbs. Is there any equipment available that could help us?
7	The reception area seems different from when I last visited in January. Has anything changed?

Table 16: List of questions used in the hardware experiments. All questions require both memory retrieval and active exploration.

H Examples of Failure Cases

Figure 13 shows two failure cases due to limitations of the current VLMs' semantic and spatial understanding capabilities. The first example shows the failure to recognize objects in the images. Despite retrieving the correct images to answer the question, the VLM fails to return the correct answer in our experiment. The second example involves counting objects through multiple images. The retrieved images contain overlapping views that cover all the blue barrels in Aisle 2. However, we often fail to give the correct number of objects. These failure cases highlight the need for improved semantic and spatial understanding of approaches that use vision and language models to improve the performance on the LA-EQA benchmark.



Question: What tools are on the floor next to tool box in room 13?

Correct answer: A drilling machine, a saw or a leveling tool (any one is correct)

Model answer: A crowbar and other construction tools

Failure reason: Limited object detection/recognition



Question: How many blue barrels do we have on the shelves in Aisle 2?

Correct answer: 7

Model answer: 5

Failure reason: Limited counting capability

Figure 13: Examples of the failure cases.

I Example of Full Results

Example 2: Full example of Mind Palace Exploration for LA-EQA (part 1/2)

Question: When was the package at the front door upstairs delivered?
Ground truth answer: Before Thursday afternoon

Reasoning over LA-EQA Iteration 1

Object to search: package at the front door upstairs

Reasoning: To determine when the package was delivered, I need to look for a package at the front door upstairs, as it is the specific object mentioned in the question.

Planning over episodic world instances

Search strategy: PAST ONLY

Reasoning on search strategy: The user question is about when the package was delivered, which is a past event. Therefore, the most efficient strategy is to recall past memory to find when the package was first observed at the front door upstairs. Exploring the present environment is not necessary as it will not provide information about the delivery time.

EpisodicExploration: Time instance to retrieve: ['friday afternoon', 'thursday afternoon', 'wednesday afternoon']

EpisodicExploration: Reasoning: These times are relevant because they are the most recent past instances before the current time. Checking these instances will help identify when the package was first observed at the front door upstairs, which is crucial to determine the delivery time.

Searching for the target object in 'friday afternoon'

Planning over areas

Area list: ['main entrance', 'Living room upstairs', 'stairs', 'kitchen and dining', 'Living room downstairs']

Area probability: [0.85, 0.7, 0.5, 0.3, 0.2]

Reasoning: The task is to locate a package at the front door upstairs and determine when it was delivered. Given the room names and their descriptions, the main entrance is the most likely location for a package delivery, as it is typically where packages are left. The living room upstairs is also a potential location since it is upstairs and could be near the front door. The stairs could be a transitional area where the package might be temporarily placed. The kitchen and dining area, while not directly related to package delivery, could be a place where someone might move the package after receiving it. The living room downstairs is less likely since it is not upstairs, but it could still be a place where the package was moved. Based on these considerations, I assign the highest probability to the main entrance, followed by the living room upstairs, stairs, kitchen and dining, and living room downstairs.

Planning over viewpoints

Viewpoints to explore: [198, 204, 207, 201, 213]

Reasoning: To locate the package at the front door upstairs, we should search places that mention a door and a cardboard box, as these are likely locations for package delivery. Place Node 198, 204, and 207 all mention a cardboard box and a door, making them relevant to the search. Additionally, Place Node 201 and 213 mention a door and are part of the main entrance, which is a common area for package delivery.

Image retrieval...

Image Analysis Answer: Object found!

Reasoning: The images show a package at the front door upstairs, indicating its presence. However, the images do not provide information about when it was delivered.

Full example of Mind Palace Exploration for LA-EQA (part 2/2)

```
# Searching for the target object in 'thursday afternoon'
# Planning over areas (main entrance)
# Episode to explore: thursday afternoon
# Planning over viewpoints
Image retrieval...
Image Analysis Answer: Object found!
Reasoning: The images show a package placed near the front door, indicating
that it has been delivered.

# Searching for the target object in 'wednesday afternoon'
# Planning over areas (main entrance)
# Planning over viewpoints
Image retrieval...
Image Analysis Answer: Object NOT found!

# Planning over areas (Living room upstairs)
# Planning over viewpoints
Image retrieval...
Image Analysis Answer: Object NOT found!

# Planning over areas (stairs)
# Planning over viewpoints
Image retrieval...
Image Analysis Answer: Object NOT found!

# Planning over areas (kitchen)
# Planning over viewpoints
Image retrieval...
Image Analysis Answer: Object NOT found!
Reasoning: The package was not found in Wednesday afternoon

# Reasoning over LA-EQA Iteration 2
Ready to answer: Yes
Reasoning: The package was already present on Thursday afternoon, as seen in
the images. Therefore, it was delivered before Thursday afternoon.

Answer: Before Thursday afternoon
Reasoning: The package was observed at the front door upstairs on Thursday
afternoon. It was not present on Wednesday afternoon, indicating it was
delivered sometime between Wednesday afternoon and Thursday afternoon.
```