
Embodied-RAG: General Non-parametric Embodied Memory for Retrieval and Generation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 There is no limit to how much a robot might explore and learn, but all of that
2 knowledge needs to be searchable and actionable. Within language research, re-
3 trieval augmented generation (RAG) has become the workhouse of large-scale non-
4 parametric knowledge, however existing techniques do not directly transfer to the
5 embodied domain, which is multimodal, data is highly correlated, and perception
6 requires abstraction. To address these challenges, we introduce Embodied-RAG,
7 a framework that enhances the foundational model of an embodied agent with a
8 non-parametric memory system capable of autonomously constructing hierarchical
9 knowledge for both navigation and language generation. Embodied-RAG handles
10 a full range of spatial and semantic resolutions across diverse environments and
11 query types, whether for a specific object or a holistic description of ambiance.
12 At its core, Embodied-RAG’s memory is structured as a semantic forest, storing
13 language descriptions at varying levels of detail. This hierarchical organization
14 allows the system to efficiently generate context-sensitive outputs across different
15 robotic platforms. We demonstrate that Embodied-RAG effectively bridges RAG
16 to the robotics domain, successfully handling over 200 explanation and naviga-
17 tion queries across 19 environments, highlighting its promise for general-purpose
18 non-parametric system for embodied agents.

19 1 Introduction

20 Humans excel as generalist embodied agents in part due to our ability to build, abstract, and reason
21 over rich memories. In contrast, current embodied agents[Chaplot et al., 2020, Khanna et al., 2024,
22 Krantz et al., 2022, Zhou et al., 2023] lack such versatile memory capabilities, limiting their ability
23 to operate effectively in unbounded and complex real-world environments. While existing methods
24 such as semantic mapping[Chaplot et al., 2020, Khanna et al., 2024] and scene graphs[Li et al., 2022,
25 Rana et al., 2023] attempt to capture spatial and contextual relationships, they largely fall short of the
26 dynamic and flexible memory, retrieval, and generative abilities exhibited by humans.

27 In the language domain, foundation models combined with non-parametric memory mechanisms
28 have achieved near human-level performance across various tasks. Retrieval-Augmented Generation
29 (RAG) [Asai et al., 2023, Chen et al., 2023, Lewis et al., 2021] has been widely adopted in the field
30 of Natural Language Processing (NLP) as a non-parametric memory mechanism over large document
31 corpora, enhancing the accuracy and relevance of responses generated by Large Language Models
32 (LLMs). Similarly, the continuous stream of experiences gathered by embodied agents forms vast
33 databases that exceed the context window limitations of LLMs.

34 However, applying RAG to embodied scenarios presents unique challenges due to key differences
35 between textual data and embodied experiences. First, while RAG relies on existing documents,
36 building memory from embodied experiences is itself a core research challenge. Current methods,

37 such as dense point clouds or scene graphs, fail to capture the full range of experiences beyond
38 object-level attributes, without relying on human-engineered schemas or exceeding memory budgets.
39 Second, unlike documents, embodied experiences have inherent correlated structure — semantically
40 similar objects are often spatially correlated and hierarchically organized so embodied experiences
41 should not be treated as independent samples. Finally, embodied observations vary in granularity
42 and structure: outdoor scenes might be sparse, while indoor environments are cluttered, and repeated
43 objects across frames can confuse LLMs, complicating retrieval.

44 We present **Embodied-RAG**, a system with two key components: Memory Construction and Retrieval
45 and Generation. In Memory Construction, Embodied-RAG autonomously builds a topological map
46 and a hierarchical semantic forest that organizes observations based on spatial correlations. This forest
47 allows retrieval at different abstraction levels (explicit, implicit, global) by matching the query with
48 corresponding memory resolution (local, intermediate, global). To mitigate perceptual hallucinations,
49 the Retrieval and Generation process incorporates parallel tree traversals scored by a language model,
50 using retrieved results for explanations or navigational actions via an LLM.

51 We evaluated Embodied-RAG, with a novel benchmark with over 200 tasks requiring multimodal
52 outputs and reasoning. Compared to Semantic Match and vanilla RAG, Embodied-RAG demonstrated
53 superior performance: (1) more robust against object detection errors on explicit queries, leveraging
54 spatial relevancy; (2) improved reasoning on implicit queries, with a 220% improvement over
55 Semantic Match and 30% over RAG; (3) better global summarization and trend analysis, where
56 Semantic Match and RAG showed poor results.

57 The key contributions and implications of this paper include:

- 58 • **Method** We introduce the system of Embodied-RAG. This method addresses problems of naively
59 apply non-parametric memories like RAG to embodied setting.
- 60 • **Task** We introduce the general task of *Embodied-RAG benchmark*, formulating semantic naviga-
61 tion and question answering under a single paradigm (Figure 2).
- 62 • **Implications** Our results and discussion provide a basis for rethinking approaches to generalist
63 robot agents based on non-parameteric memories.

64 2 Task

65 We introduce the Embodied-RAG benchmark, which contains queries from the cross-product of
66 {explicit, implicit, global} questions with potential {navigational action, language} generation outputs.

67 A task consists of:

- 68 • **Query:** The content can be explicit (e.g. a particular object instance), implicit (e.g. looking for
69 adequacy, instruction with more pragmatic understanding required), or global. The request might
70 pertain to a location or general vibe.
- 71 • **Experience:** The experience is a sequence of egocentric visual perception and odometry, occurring
72 in indoor, outdoor, or mixed environments.
- 73 • **Output:** The expected output can be both navigation actions with language descriptions (Fig 2
74 top, Fig. 1 c-1), or language explanations (Fig 2 bottom, Fig. 1 c-2).

75 Example tasks are shown in Fig. 2, with instances of explicit, implicit, and global queries in Fig. 1.

76 3 Method

77 3.1 Memory Construction

78 The memory construction process of Embodied-RAG consists of two parts: a topological map and a
79 semantic forest.

81 **Topological map** We employ a topological graph composed of nodes with the following attributes:

- 82 • **Position Information:** Allocentric coordinates (x, y, z) and the yaw angle θ .
- 83 • **Image Path:** Each node contains a path to an associated ego-centric image.
- 84 • **Captions:** Generated by a vision-language model, these captions provide object-level natural
85 language textual descriptions of the image.

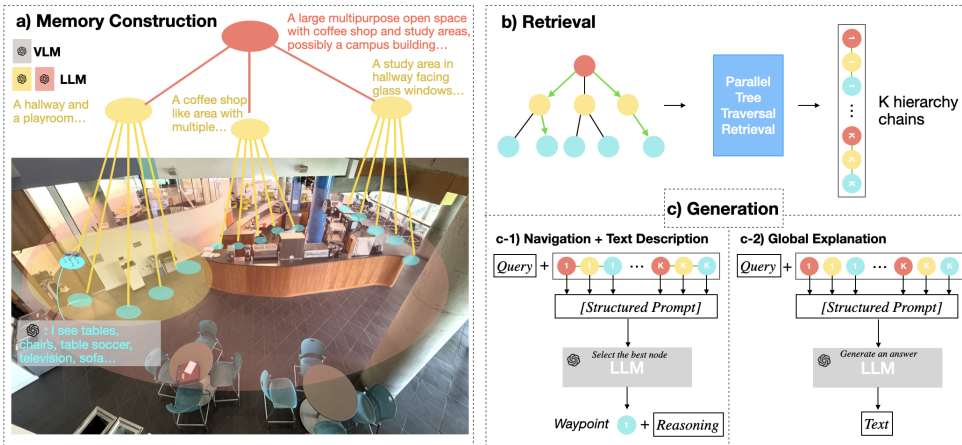


Figure 1: **Embodied-RAG method overview.** (a) Memory is constructed by hierarchically organizing the nodes of the topological map into a semantic forest. (b) The memory in (a) can be retrieved for a query, with parallelized tree traversals. (c) Navigation actions with text outputs, or global explanations can be generated for the query, with using the retrieval results as LLM contexts.

86 The nodes form a topological map (blue nodes in Fig. 1), eliminating the need for specific control
 87 parameters like velocity and yaw, which often vary across different drive systems. This abstraction
 88 enables compatibility with any local planner, regardless of the robot’s embodiment. Furthermore,
 89 the topological structure is far more memory-efficient than traditional metric maps [Chaplot et al.,
 90 2020, Min et al., 2021, Shafullah et al., 2022], allowing for efficient scaling in both large outdoor
 91 and complex indoor environments. Our experiments show that this approach successfully navigates
 92 kilometer-scale simulated environments.

93 **Semantic Forest** We use a separate tree structure, referred to as a semantic forest, to capture meaning
 94 at various spatial resolutions. The nodes of this forest are those of the topological map, with the
 95 non-leaf nodes capturing larger spaces at a thinner density of semantic specificity. First, we create
 96 the forest through hierarchical clustering. Since spatially approximate leaf nodes exhibit semantic
 97 correlations, we employ an agglomerative clustering mechanism [Sneath and Sokal, 1973] to group
 98 nodes based on their physical positions assigning the mean position of the leaves.

99 This iterative process continues until a root node is formed, stopping when no further relevance is
 100 found based on a threshold set by the algorithm. Once we have a complete forest with one or more
 101 root nodes, each non-leaf node receives a language description. We achieve this by prompting a large
 102 language model (LLM, e.g., GPT-4) to generate a abstraction that encompasses the descriptions of its
 103 direct child nodes (see website for the prompting). This process is conducted bottom-up, starting
 104 from the leaf nodes and moving up to the parent nodes. We parallelized this process across all nodes
 105 at the same hierarchy.

106 3.2 Retrieval

107 We run the following *process*, which takes a single tree as input and outputs a single leaf node.
 108 Starting by visiting the root node, we run BFS with LLM selection; we ask *LLM_Selector* to choose
 109 the best child node of the currently visited node based on compatibility with the given query. For
 110 example, if the query is “find me a place that is bright and quiet but has some presence of people,”
 111 we prompt the LLM to select the best description among the children of the currently visited node.
 112 We then visit the selected best child node and iterate this process until we reach a leaf node. Once
 113 we obtain k leaf nodes ($\frac{k}{N}$ nodes from each tree) by running this process $\frac{k}{N}$ times for each of the N
 114 trees, we obtain the “chain” from the selected node to the root node. The $\frac{k}{N}$ processes are parallelized
 115 across the N trees. The set of these best k chains is the retrieval output, containing semantics at
 116 all scales for any specific location that corresponds to the leaf scale. Embodied-RAG unifies the
 117 retrieval process to handle explicit, implicit, and global queries, producing both explanations and
 118 navigational actions as outputs. Note, these hierarchies and corresponding trees allow for querying
 119 automatically created semantic regions, something particularly useful for outdoor navigation where
 120 walls and structures cannot be used to determine function.

121 3.3 Generation

122 We pass the retrieved k best chains as part of a context, for the LLM to generate navigation and text
123 description (Fig. 2 top) or global explanations (Fig. 2 bottom). Given the query and the k chains,
124 we prompt the LLM to “select” a waypoint with a reasoning, or to “explain” (prompt in our project
125 website).

126 **Navigation** We select a waypoint (a leaf node of the semantic forest) and use a planner to generate
127 navigational actions—sequences of (torque, velocity) pairs— to reach the waypoint. To select this
128 waypoint, we ask the LLM to choose the best single leaf node, together with textual reasoning, using
129 the query and the chain as input. Again, including the entire chain as input ensures that a waypoint
130 can be generated for implicit navigation tasks as well.

131 **Text Answers** As depicted in Figure 1 (c), we concatenate the k chains as part of the prompt to the
132 LLM. We ask it to generate an answer to the query based on the k retrieved chains. The spatial scale
133 of attention in each node of the chain facilitate the LLM to generate responses at any semantic scale
134 (explicit, implicit, general) based on the retrieved result.

135 4 Results

Table 1: Comparison of Methods on different Embodied-RAG Benchmarks.

Env.	Explicit			Implicit			Global		
	Embodied-RAG	RAG	Sem.	Embodied-RAG	RAG	Sem.	Embodied-RAG	RAG	Sem.
Small	0.955	0.955	0.955	1.000	0.818	0.364	4.88	3.67	-
Large	0.977	0.947	0.895	0.914	0.695	0.426	4.86	2.43	-
Total	0.969	0.949	0.877	0.926	0.706	0.410	4.87	2.68	-

136 We present **quantitative** results in Table 1, demonstrating the effectiveness of Embodied-RAG across
137 explicit, implicit, and global retrieval tasks. We also classify environments as small or large based
138 on the number of mapped nodes. Embodied-RAG consistently outperforms RAG and Semantic
139 Match across all tasks and environments. While all methods perform well on explicit queries,
140 Embodied-RAG provides a slight advantage due to its hierarchical structure. For implicit queries,
141 RAG and Semantic Match performance drops significantly, especially in large environments, while
142 Embodied-RAG remains robust. On global questions, Embodied-RAG excels, while Semantic Match,
143 lacking summarization and reasoning, cannot be applied.

144 We conducted a **qualitative** comparison between Embodied-RAG and baseline models.

145 *Implicit Query:* Where can I buy drinks? Embodied-RAG correctly identifies a food service area,
146 while the baselines return irrelevant answers like a refrigerator or water fountain. These results reflect
147 a mismatch between the user’s intent (to buy) and the retrieved objects. Embodied-RAG performs
148 multi-step reasoning and retrieves more suitable locations, such as counters or vending machines,
149 matching the user’s intent.

150 *Global Query:* As illustrated in Figure 2, Embodied-RAG accurately describes the environment as
151 a suburban neighborhood with a park, using its hierarchical structure to provide a cohesive view.
152 In contrast, RAG retrieves isolated nodes, resulting in a fragmented and redundant interpretation.
153 Embodied-RAG integrates elements like trees and shrubs into the broader park context, offering a
154 more human-like understanding.

155 5 Conclusion

156 We present Embodied-RAG, a system that captures spatial memory at any resolution and generates
157 responses for both navigation and explanation requests. We also introduce the Embodied-RAG
158 benchmark, unifying semantic navigation and question answering. Our results show that Embodied-
159 RAG handles implicit and global queries, as well as ambiguous human requests, demonstrating its
160 potential for integrating large non-parametric memories into robotics models. We look forward to
161 future extensions involving manipulation and dynamic environments, enabling robotics tasks beyond
162 current memory-constrained methods.

163 **References**

- 164 Akari Asai, Sewon Min, Zexuan Zhong, and Danqi Chen. Acl 2023 tutorial: Retrieval-based language
165 models and applications. *ACL 2023*, 2023.
- 166 Chaplot, Russ R, et al. Object goal navigation using goal-oriented semantic exploration. *NeurIPS*,
167 33, 2020.
- 168 Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. Benchmarking large language models in
169 retrieval-augmented generation, 2023.
- 170 Khanna, Roozbeh, et al. Goat-bench: A benchmark for multi-modal lifelong navigation.
171 *arXiv:2404.06609*, 2024.
- 172 Jacob Krantz, Stefan Lee, Jitendra Malik, Dhruv Batra, and Devendra Singh Chaplot. Instance-specific
173 image goal navigation: Training embodied agents to find object instances. *CVPR*, 2022.
- 174 Patrick Lewis, Douwe Kiela, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks,
175 2021.
- 176 Li, Fuchun, et al. Embodied semantic scene graph generation. In Aleksandra Faust, David Hsu, and
177 Gerhard Neumann, editors, *CoRL*. PMLR, 2022.
- 178 So Yeon Min, Devendra Singh Chaplot, Pradeep Ravikumar, Yonatan Bisk, and Ruslan Salakhutdinov.
179 Film: Following instructions in language with modular methods. *ICLR*, 2021.
- 180 Krishan Rana et al. Sayplan: Grounding large language models using 3d scene graphs for scalable
181 task planning. In *CoRL*, 2023.
- 182 Nur Muhammad Mahi Shafiullah, Chris Paxton, Lerrel Pinto, Soumith Chintala, and Arthur Szlam.
183 Clip-fields: Weakly supervised semantic fields for robotic memory. *arXiv: Arxiv-2210.05663*,
184 2022.
- 185 Peter HA Sneath and Robert R Sokal. *Numerical Taxonomy: The Principles and Practice of*
186 *Numerical Classification*. W.H. Freeman, 1973.
- 187 Gengze Zhou, Yicong Hong, and Qi Wu. Navgpt: Explicit reasoning in vision-and-language
188 navigation with large language models. *arXiv:2305.16986*, 2023.

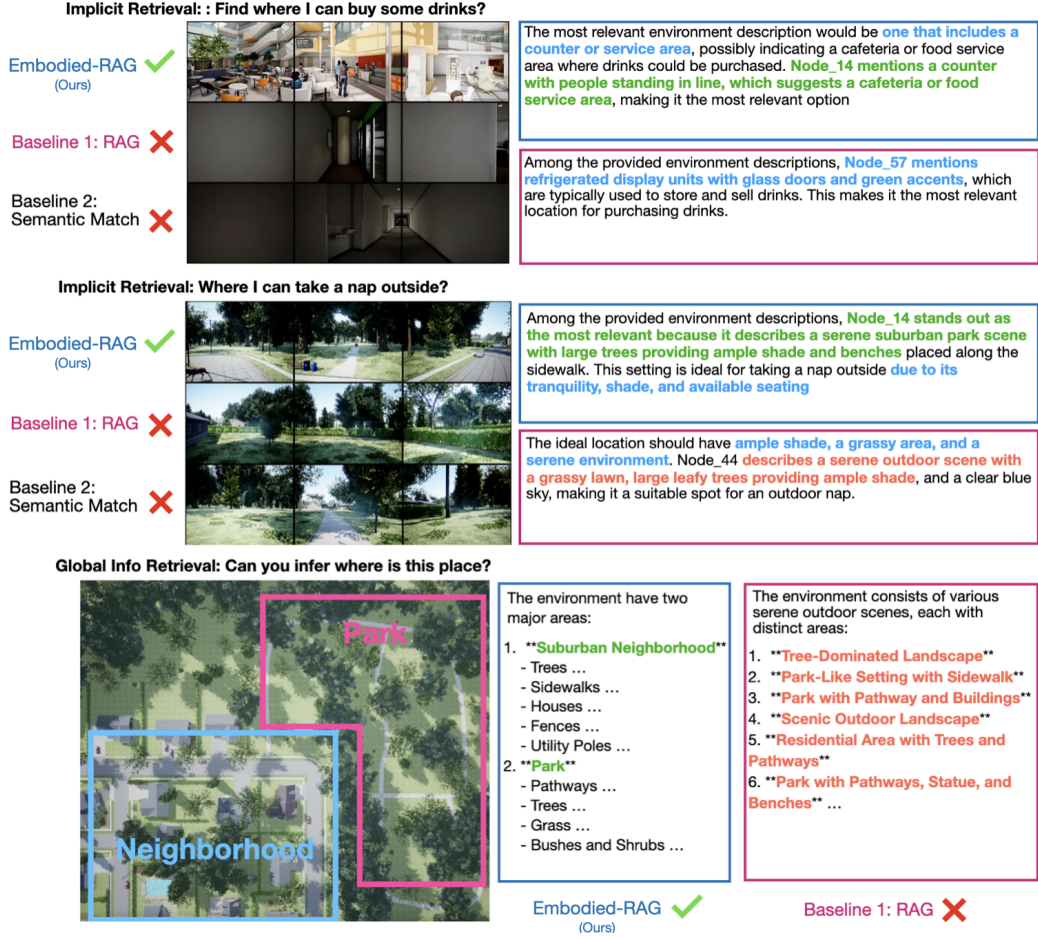


Figure 2: Example reasoning of Embodied-RAG and RAG for generation tasks are highlighted in blue and pink boxes, respectively.

189 A Computational Efficiency

190 Both memory construction and retrieval have a computational complexity of $O(\log N)$, where N
 191 represents the number of nodes in the environment. This choice allows us to efficiently scale to
 192 larger environments, as the time complexity only increases logarithmically with the number of nodes.
 193 Additionally, when performing the k retrievals, we execute them in parallel to minimize the overall
 194 time cost. In our real-life experiments, the time costs are demonstrated in the supplementary video,
 195 which is 8x fast-forwarded. On average, a single retrieval takes around 20 seconds in most of our
 196 environments, and the travel time depends on the speed of the specific embodiment in use.

197 B Ablation

198 We investigate the impact of $k \in \{1, \text{GPT4 Token Limit}\}$ on Embodied-RAG and RAG in Figure
 199 3. A total of 15 experiments were conducted for each k in each environment. We observe that with
 200 larger k , both RAG and Embodied-RAG show improved performance, but this improvement plateaus
 201 at higher values. RAG still fails to capture the larger holistic resolution with just more object-level
 202 nodes and cannot adequately solve the implicit/general queries, further justifying our hierarchy and
 203 tree selection approach.

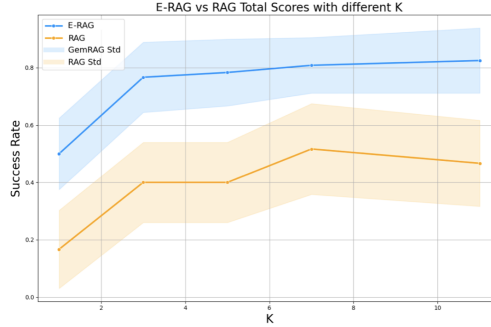


Figure 3: Effect of total number of K searches or K retrievals

204 **C Limitations and Future work**

205 We primarily focused on semantic forests rather than a topological map. Therefore, we may not be
 206 robust in obstacle avoidance involving dynamic objects and people. Furthermore, Embodied-RAG
 207 currently struggles with requests that require precise counting of objects at a small scale (e.g., “How
 208 many chairs are there around the red table?”). This limitation arises because the agglomerative
 209 clustering of the semantic forest does not consider multi-view consistency. Future work could
 210 incorporate multi-view consistency in the hierarchies of the semantic forest with a learned or pre-
 211 trained mechanism to cluster with positional information (e.g. utilizing a LLM).