IMPROVING MULTI-STEP RAG WITH HYPERGRAPH-BASED MEMORY

Anonymous authors

000

001

002003004

006

008 009

010 011

012

013

014

015

016

018

019

020

021

022

024

025

026

027

028

029

031 032 033

034

037

038

040

041 042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Multi-step retrieval-augmented generation (RAG) has become a widely adopted strategy for enhancing large language models (LLMs) on tasks that demand global comprehension and intensive reasoning. Many RAG systems incorporate a working memory module to consolidate retrieved information. However, existing memory designs function primarily as passive storage that accumulates isolated facts for the purpose of condensing the lengthy inputs and generating new sub-queries through deduction. This static nature overlooks the crucial high-order correlations among primitive facts, the compositions of which can often provide stronger guidance for subsequent steps. Therefore, their representational strength and impact on multi-step reasoning and knowledge evolution are limited; resulting in fragmented reasoning and weak global sense-making capacity in extended contexts. We introduce HGMEM, a hypergraph-based memory mechanism that extends the concept of memory beyond simple storage into a dynamic, expressive structure for complex reasoning and global understanding. In our approach, memory is represented as a hypergraph whose hyperedges correspond to distinct memory units, enabling the progressive formation of higher-order interactions within memory. This mechanism connects facts and thoughts around the focal problem, evolving into an integrated and situated knowledge structure that provides strong propositions for deeper reasoning in subsequent steps. We evaluate HGMEM on several challenging datasets designed for global sense-making. Extensive experiments and in-depth analyses show that our method consistently improves multi-step RAG and substantially outperforms strong baseline systems across diverse tasks.

1 Introduction

Large Language Models (LLMs) have achieved impressive progress across a wide range of tasks. Nevertheless, their ability to handle long-context problems remains fundamentally constrained by finite context windows, making it impractical to simply pack all relevant information into the input. This limitation has spurred the development of retrieval-augmented generation (RAG), which supplements LLMs with external knowledge retrieval. Yet, single-step retrieval often proves insufficient for complex queries, motivating the shift toward multi-step RAG methods that iteratively interleave retrieval with reasoning.

To effectively capture dependencies across steps and condense the lengthy processing history, many recent approaches incorporate working memory mechanisms inspired by human cognition (Lee et al., 2024; Zhong et al., 2024). A straightforward solution is to let LLMs summarize the interaction history and environment feedback into a textual description of the current problem-solving state. This strategy has been widely adopted since early work on multi-step RAG (Li et al., 2023; Trivedi et al., 2023) as well as in commercial systems (Jones, 2025; Shen & Yang, 2025). However, such unstructured memory mechanisms cannot be manipulated with sufficient accuracy across steps and often lose the ability to back-trace references to retrieved texts. Consequently, recent research has shifted toward structured or semi-structured working memory, typically with predefined schemas such as relational tables (Lu et al., 2023), knowledge graphs (Oguz et al., 2022; Xu et al., 2025), or relational/event-centric bullet points (Wang et al., 2025).

Existing memory mechanisms often treat memory as static storage that continually accumulates meaningful but primitive facts. This view overlooks the evolving nature of human working mem-

ory, which incrementally incorporates higher-order correlations from previously memorized content. Such capacity is crucial for solving complex problems that involve long-context evidence, *e.g.*, global sense-making tasks. In these scenarios, the required evidence is often composed of complex knowledge structures that extend beyond predefined schemas, and reasoning over long lists of primitive facts is both inefficient and prone to confusion with mixed or irrelevant information. Current memory mechanisms in multi-step RAG systems lack these abilities, preventing memory from effectively guiding LLMs' interaction with external data sources. These limitations highlight the need for memory with stronger representational capacity.

In this paper, we propose a hypergraph-based memory mechanism (HGMEM) for multi-step RAG systems, which enables memory to evolve into more expressive structures that support complex reasoning and global understanding. Hypergraphs, as a generalization of graphs, are particularly well-suited for this purpose (Feng et al., 2019). In our design, memory is structured as a hypergraph composed of hyperedges, each treated as a distinct memory point that represents a specific perspective of the memorized information. Initially, these memory points encode low-order primitive facts. As the LLM interacts with external environments, higher-order correlations among memory points gradually emerge and are progressively integrated into the memory through update, insertion, and merging operations. At each step prior to response generation, the LLM examines the current memory and adaptively generates subqueries, enabling both focused local investigation and broad global exploration.

This rich and structured memory facilitates broader contextual awareness and stronger reasoning in real-world applications by offering several advantages. First, it maintains an **integrated body of knowledge** around the focal problem by synthesizing primitive evidence and intermediate thoughts, typically going *beyond predefined schemas* and providing a *global perspective* over the evidence. Second, it offers **structured and accurate guidance** for the LLM's sustained interactions in two ways: (1) enabling subsequent reasoning to start from representational propositions rather than from a long list of disparate primitive facts; and (2) guiding subquery generation and retrieval in a more accurate manner, leveraging each hyperedge without mixing irrelevant information.

We conduct extensive experiments on several challenging tasks involving global sense-making questions within long context. The results show that our HGMEM achieves significant improvements over competitive RAG baselines, confirming the advantages.

2 Related Work

2.1 WORKING MEMORY MECHANISMS FOR MULTI-STEP RAG

Starting from ReAct (Yao et al., 2023), many multi-step RAG systems have incorporated reflections to integrate available information for subsequent decisions. These reflections can be regarded as a simple form of memory. With the development of structured indexing for RAG, working memory also borrows this idea. Prevailing studies (Li et al., 2023; Shen & Yang, 2025; Chhikara et al., 2025; Xu et al., 2025) save agent behavior, such as task decomposing, execution tracking, and result verification, to manage task context more effectively, representing a step toward explicit working memory for complex multi-agent coordination. This idea also matured in chain-of-thought (CoT) and multiround RAG, where working memory is represented as iteratively updated records of reasoning steps or retrieved evidence. For example, IRCOT (Trivedi et al., 2023) and ComoRAG (Wang et al., 2025) employ a dynamic memory workspace to iteratively consolidate past knowledge or steps and incorporate new evidence, supporting scalable and iterative reasoning across multiple steps.

2.2 RAG WITH STRUCTURED INDEX

There is a long line of work that studies managing extended corpora through structured indexing for RAG. Though different from our focus on working memory mechanism, these work can be viewed as building structured (and static) long-term memory thus is relevant. Specifically, tree-structured methods, such as Sarthi et al. (2024), Fatehkia et al. (2024), and Tao et al. (2025), organize text chunks or entity hierarchies, enabling multi-level or bidirectional retrieval to enhance context integration. Graph-based methods, such as GraphRAG (Edge et al., 2024) and LightRAG (Guo et al., 2024), construct entity graphs and community summaries or utilize graph-enhanced indexing for dual-level retrieval, thereby improving global reasoning, efficiency, and response diversity.

In addition, to address the finite context window of LLMs, a range of memory mechanisms has been proposed that simulate long contexts or dialog histories as long-term memory to improve multi-step RAG systems. According to memory representation, most methods (Gutierrez et al., 2024; Lee et al., 2024; Li et al., 2024; Gutiérrez et al., 2025) adopt contextual memory that explicitly stores retrieved information, complementing LLMs' parameters, which directly organize long-contexts or passages into pre-defined structures such as graphs, trees, or hierarchical summaries, and retrieve relevant knowledge for multi-step reasoning. For instance, Chen et al. (2023); Gutierrez et al. (2024); Li et al. (2024) construct explicit knowledge graphs or passage indices, and then integrate information via multi-hop retrieval or agent-guided graph exploration. Gutiérrez et al. (2025); Rasmussen et al. (2025) enhance retrieval with deeper passage integration. In contrast, another line of work employs parametric memory by utilizing LLMs or dedicated modules to generate auxiliary parametric representations to enhance RAG execution (Qian et al., 2025).

However, these studies merely leverage their memories as basic storages that accumulate primitive facts, ignoring the underlying correlations among memorized content for global sense-making tasks. Therefore, designing a flexible and generalized memory mechanism that supports modeling high-order correlations is crucial for better guiding LLMs to conduct multi-step RAG.

3 METHODOLOGY

 We introduce HGMEM, the hypergraph-based memory mechanism designed to facilitate better contextual awareness and reasoning in multi-step RAG settings with structured data sources, especially for long-context tasks that require complex global sense-making.

3.1 PROBLEM FORMULATION

In this work, we consider the kind of tasks for LLMs to resolve a query based on a given document. Besides the plain texts, we assume that the document has been preprocessed into a graph through an offline graph-building stage, where entities and relationships are extracted from the document passage. Formally, let us denote the document as \mathcal{D} segmented into a set of small manageable text chunks $\{d_1, d_2, ..., d_{|\mathcal{D}|}\}$, and the derived graph as \mathcal{G} composed of nodes $\mathcal{V}_{\mathcal{G}}$ and edges $\mathcal{E}_{\mathcal{G}}$ corresponding to the extracted entities and relationships, respectively. Each node $v \in \mathcal{V}_{\mathcal{G}}$ or edge $e \in \mathcal{E}_{\mathcal{G}}$ is associated with the source text chunks in which its embodied entity/relationship appears, which is recorded during the offline graph construction. Meanwhile, the nodes, edges and text chunks are embedded into high-dimensional vectors for vector-based retrieval. For resolving the query, LLMs have access to both the document and its derived graph as structured data sources.

3.2 MULTI-STEP RAG SYSTEM WITH MEMORY

When dealing with tasks requiring comprehensive understanding, especially over long context, RAG systems usually resort to multi-step approaches with an underlying memory mechanism, where retrieval operations are interleaved with intermediate reasoning to support broader contextual awareness.

Given a target query \hat{q} , the LLM iteratively interacts with \mathcal{D} and \mathcal{G} while managing a memory \mathcal{M} to store relevant information for ultimately resolving \hat{q} . During each interaction step t, the LLM judges whether the content of current memory has been sufficient with respect to the target query. If the model judges the memory has contained sufficient information, it immediately produces a response. Otherwise, it analyzes current memory and generates several subqueries $\mathcal{Q}^{(t)}$ that aim at fetching more information from D and G to enrich the memory.

As illustrated in Figure 1 (i), at the t-th step, if the LLM proceeds to generate subqueries $\mathcal{Q}^{(t)}$ based on current memory $\mathcal{M}^{(t)}$ maintained until the previous step, it retrieves a set of the most relevant entities $\mathcal{R}_v(\mathcal{Q}^{(t)})$ from \mathcal{G} using vector-based matching. Here, $\mathcal{R}_v(\mathcal{Q})$ defines the operation of retrieving entities with the highest relevance to a query set \mathcal{Q} :

$$\mathcal{R}_{v}(\mathcal{Q}) = \bigcup_{q \in \mathcal{Q}} \underset{v \in \mathcal{V}_{\mathcal{G}}}{\operatorname{argmax}}(\operatorname{sim}(\mathbf{h}_{q}, \mathbf{h}_{v})), \tag{1}$$

¹In this work, we just use the graph construction tool open-sourced by LightRAG (Guo et al., 2024).

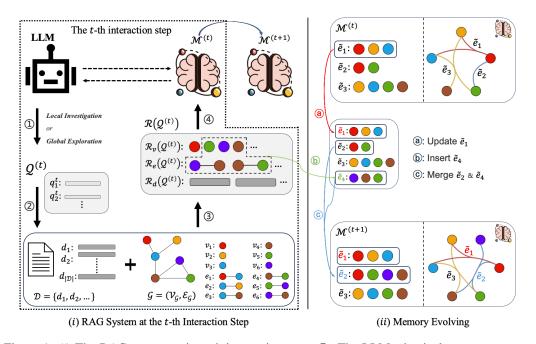


Figure 1: (i) The RAG system at its t-th interaction step. ①: The LLM adaptively generates a set of subqueries $\mathcal{Q}^{(t)}$ for either local investigation or global exploration (see Section 3.4). ②: $\mathcal{Q}^{(t)}$ are used to retrieve information from \mathcal{D} and \mathcal{G} . ③: $\mathcal{R}(\mathcal{Q}^{(t)})$ is obtained through graph-based indexing and vector-based matching. ④: The LLM analyzes $\mathcal{R}(\mathcal{Q}^{(t)})$ and evolves current memory $\mathcal{M}^{(t)}$ into $\mathcal{M}^{(t+1)}$. (ii) The structure of our proposed hypergraph-based memory that evolves through update, insertion and merging operations.

where n_v is the limit number of retrieved entities per query, \mathbf{h}_q is the vector representation of q, \mathbf{h}_v is the vector representation of v, and $\mathrm{sim}(\cdot,\cdot)$ is the cosine similarity function. Then, via graph-based indexing, the relationships and text chunks associated with these entities are also obtained, denoted as $\mathcal{R}_e(\mathcal{Q}^{(t)})$ and $\mathcal{R}_d(\mathcal{Q}^{(t)})$, respectively. Subsequently, the LLM analyzes and consolidates these information into the memory, evolving memory into $\mathcal{M}^{(t+1)}$, which can be formalized as

$$\mathcal{M}^{(t+1)} \leftarrow \text{LLM}(\mathcal{R}(\mathcal{Q}^{(t)}), \mathcal{M}^{(t)}),$$
 (2)

$$\mathcal{R}(\mathcal{Q}^{(t)}) = (\mathcal{R}_v(\mathcal{Q}^{(t)}), \mathcal{R}_e(\mathcal{Q}^{(t)}), \mathcal{R}_d(\mathcal{Q}^{(t)})). \tag{3}$$

Note that, at the initial step (t=0), we treat the target query \hat{q} as a special subquery belonging to $\mathcal{Q}^{(0)}$, *i.e.* $\mathcal{Q}^{(0)} = \{\hat{q}\}$. Further details about the memory storage, subquery generation and the dynamic of memory evolving will be elaborated in Section 3.3, Section 3.4 and Section 3.5, respectively.

3.3 Hypergraph-based Memory Storage

When the LLM interacts with the document \mathcal{D} and the graph \mathcal{G} , it continuously consolidates relevant information into the memory storage. In our RAG system, the entire memory is modeled as a hypergraph $\mathcal{M}=(\tilde{\mathcal{V}}_{\mathcal{M}},\tilde{\mathcal{E}}_{\mathcal{M}})$, where $\tilde{\mathcal{V}}_{\mathcal{M}}=\{\tilde{v}_1,\tilde{v}_2,...\}$ is the node set and $\tilde{\mathcal{E}}_{\mathcal{M}}=\{\tilde{e}_1,\tilde{e}_2,...\}$ is the hyperedge set. Similar to those nodes in \mathcal{G} , every node $\tilde{v}\in\tilde{\mathcal{V}}_{\mathcal{M}}$ in \mathcal{M} also has its associated text chunks, along with the name and description of its embodied entity. Meanwhile, it is noteworthy that the set of nodes $\tilde{\mathcal{V}}_{\mathcal{M}}$ in the memory correspond to the entities dynamically identified over past interactions, which are not equivalent to those entities statically existing in \mathcal{G} . As for the hyperedges, unlike those ordinary edges $\mathcal{E}_{\mathcal{G}}$ in \mathcal{G} that connect at most two nodes, each hyperedge $\tilde{e}\in\tilde{\mathcal{E}}_{\mathcal{M}}$ in \mathcal{M} can connect arbitrary number (two or more) of nodes, thus capable of modeling high-order correlation among different nodes. Particularly, the hyperedges can be treated as separate memory points,

 $^{^2}$ We also use vector-based filtering to keep at most n_e relationships and n_d text chunks during graph-based indexing.



Figure 2: An illustration of memory evolving dynamics. Each point is equivalent to a hyperedge in the memory hypergraph. The memory evolves from $\mathcal{M}^{(t)}$ into $\mathcal{M}^{(t+1)}$ through update, insertion and merging operations.

each of which can represent a certain aspect of the entire information stored in current memory, as shown in Figure 1 (ii).

As a result, the whole memory as a hypergraph can flexibly support complex relational modeling, ensuring strong expressiveness closely-related to the intrinsic structure of external data sources.

3.4 Adaptive Memory-based Subquery Generation

As described in Section 3.2, at each step t of our RAG workflow, with respect to the target query, the LLM determines whether to immediately produce a response or proceed to acquire more information from \mathcal{D} and \mathcal{G} . If current memory $\mathcal{M}^{(t)} = (\tilde{\mathcal{V}}_{\mathcal{M}}^{(t)}, \tilde{\mathcal{E}}_{\mathcal{M}}^{(t)})$ is deemed insufficient, the LLM first analyzes $\mathcal{M}^{(t)}$ and generate several subqueries $\mathcal{Q}^{(t)}$ indicating what to further explore. Specifically, we design an adaptive memory-based subquery generation strategy that produces subqueries for either local investigation or global exploration:

- (i) Local Investigation: When there are some specific memory points (a.k.a. hyperedges) worth more in-depth investigation, The LLM gives a set of subqueries, each of which targets at investigating a specific memory point more deeply. Suppose a query $q \in \mathcal{Q}^{(t)}$ is especially pertinent to a memory point $\tilde{e}_q \in \tilde{\mathcal{E}}_{\mathcal{M}}^{(t)}$, the retrieval operation triggered by this query q would be guided by all the entities belonging to this memory point \tilde{e}_q , i.e. $\{\tilde{v} | \tilde{v} \in \tilde{e}_q\}$.
- (ii) Global Exploration: When there are unexplored aspects transcending the scope of current memory, which necessitates broader exploration, the LLM resorts to generating subqueries subsequently used to explore more information from the whole external data sources, not pertinent to any existing memory point.

Under such strategy, the RAG system is able to adaptively combine both local investigation and global exploration for more flexible information retrieval during interaction with external data sources. More details and related prompts are given in Appendix C.

3.5 DYNAMIC OF MEMORY EVOLVING

Once a set of subqueries $\mathcal{Q}^{(t)}$ have been generated at the t-th step, following Equation 2, the LLM analyzes the retrieved information $\mathcal{R}(\mathcal{Q}^{(t)})$ and consolidates useful content into current memory $\mathcal{M}^{(t)}$, resulting in the evolved memory $\mathcal{M}^{(t+1)}$. As shown in Figure 1 (ii), on the basis of hypergraph-based memory storage, the dynamic of memory evolving in our proposed HGMEM involves the following three types of operations:

- *Update*. According to the retrieved information, if there are certain existing memory points whose descriptions should be modified, the update operation will revise the descriptions of corresponding hyperedges without changing their subordinate entities.
- *Insertion*. The insertion operation should be evoked when some content of the retrieved information are suitable to be inserted as additional memory points into current memory, which creates new hyperedges into the hypergraph.
- Merging. After insertion and update, the LLM inspects current memory and selectively merge
 existing memory points that are more suitable to constitute a single semantically/logically cohesive

unit. This operation builds further associations among different memory points, facilitating the resolving of queries that require complicated sense-making with disparate facts.

In this way, besides continuously accumulating primitive facts during the LLM's interactions with external data sources, the memory also gradually evolves into more sophisticated forms capturing higher-order correlations for complex relational modeling. Figure 2 gives a concrete example illustrating the dynamic of memory evolving.

3.6 Memory-enhanced Response Generation

When the LLM exceeds its maximum interaction steps or the content in current memory $\mathcal{M}^{(t)} = (\tilde{\mathcal{V}}_{\mathcal{M}}^{(t)}, \tilde{\mathcal{E}}_{\mathcal{M}}^{(t)})$ has been deemed sufficient, a response is immediately produced according to the information stored in current memory. Concretely, besides the descriptions of all memory points (i.e. $\tilde{\mathcal{E}}_{\mathcal{M}}^{(t)}$), a set of text chunks associated with all the entities $\tilde{\mathcal{V}}_{\mathcal{M}}^{(t)}$ contained in the memory are also provided to the LLM for producing the final response.

4 EXPERIMENTAL SETTINGS

4.1 DATASETS

We use generative sense-making question answering (QA) and long narrative understanding tasks for evaluation. For generative sense-making QA, similar to the setups used in previous works (Edge et al., 2024; Guo et al., 2024), we retain a portion of long documents with more than 100k tokens from **Longbench V2** (Bai et al., 2025). From each retained document, we use GPT-40 to generate several global sense-making queries that satisfy the following requirements: 1) The queries should target at the overall understanding of the whole provided documents, instead of only concentrating on several specific phrases or sentence pieces. 2) The queries should require high-level understandings and global reasoning over disparate evidences scattered across the whole paragraph. For long narrative understanding, we use three public benchmarks including **NarrativeQA** (Kociský et al., 2018), **NoCha** (Karpinska et al., 2024) and **Prelude** (Yu et al., 2025). Both tasks require global comprehension and complex sense-making over disparate evidences across long contexts. Details about the usage and statistics of data used in our experiments are given in Appendix A.

4.2 IMPLEMENTATION DETAILS

Offline Graph Construction. For all the datasets used in our experiments, we first segment every document into text chunks of 200 tokens with 50 overlapping tokens between adjacent chunks. Then, GPT-40 is utilized to preprocess each of the chunkized documents into a graph using the open-sourced tool provided by LightRAG (Guo et al., 2024). After building the graph, we adopt *bge-m3* (Chen et al., 2024) as the embedding model to convert all the entities, relationships and text chunks into vector representations managed by *nano vector database*.

System Deployment and Configuration. Our RAG system is comprised of the backbone LLM and the hypergraph-based memory. We choose GPT-40 and Qwen2.5-32B-Instruct as the representatives of advanced closed-source and open-source LLMs, respectively. During experiments, GPT-40 is accessed through official API while Qwen2.5-32B-Instruct is locally deployed with VLLM (Kwon et al., 2023). For the configuration of LLM inference, we set the temperature to 0.8 and the maximum number of output tokens to 2,048. As for the hypergraph-based memory, we employ the *hypergraph-db*³ package to maintain and manage the hypergraph at runtime. The vector representations of the nodes, hyperedges and associated text chunks in the hypergraph are also generated by *bge-m3* embedding model.

4.3 BASELINES AND EVALUATION METRICS

In our experiments, we compare our propsed HGMEM to two types of baseline methods, *i.e.* traditional RAG and multi-step RAG, which utilize plain texts and/or graph-structured data sources.

³https://github.com/iMoonLab/Hypergraph-DB

Table 1: The overall experimental results on four benchmarks. The second column "Working Memory" distinguishes whether corresponding method is equipped with a working memory that enhances LLMs during RAG execution. The best scores in each dataset are **bolded**. HGMEM consistently outperforms other comparison methods across all datasets.

Туре	Working Memory	Method	Longbench		NarrativeQA	NoCha	Prelude
турс			Comprehensiveness	Diversity	Acc (%)	Acc (%)	Acc (%)
GPT-40							
Traditional RAG	×	NaiveRAG	61.62	64.20	52.00	67.46	60.00
	×	GraphRAG	60.39	64.02	53.00	70.63	59.26
	×	LightRAG	61.55	63.37	44.00	71.43	61.48
	×	HippoRAG v2	58.92	61.27	34.00	72.22	54.81
Multi-step RAG		DeepRAG	63.62	65.98	45.00	67.46	56.30
	✓	ComoRAG	62.18	65.82	54.00	63.49	54.07
Ours		HGMEM	65.73	69.74	55.00	73.81	62.96
Qwen2.5-32B-Inst	truct						
Traditional RAG	×	NaiveRAG	61.41	62.25	37.00	64.29	52.59
	×	GraphRAG	60.78	62.16	44.00	62.70	50.37
	×	LightRAG	60.82	62.73	40.00	59.52	60.74
	×	HippoRAG v2	56.66	60.80	33.00	68.25	51.85
Multi-step RAG		DeepRAG	61.45	- 63.56 -	44.00	66.40	51.11
	✓	ComoRAG	60.74	61.28	44.00	57.60	50.37
Ōurs		HGMEM	64.18	66.51	<u>5</u> 1.00	70.63	62.22

Among these methods, DeepRAG (Guan et al., 2025) and ComoRAG (Wang et al., 2025) are equipped with a working memory while the others are not. The details of these comparison methods can be found in Appendix B. To ensure fair comparison, all baselines operate on a similar number of retrieved passages. In the case of single-step RAG, this means retrieving the same average number of text chunks as our HGMEM. For multi-step RAG methods, we approximate comparability by constraining them to rewrite the same maximum number of subqueries and perform the same maximum number of steps, while requiring retrieval of the same average number of chunks per step.

For generative sense-making QA, we adopt the following two metrics to assess the qualities of model responses: 1) **Comprehensiveness** measures how well the model response comprehensively covers and addresses all aspects and necessary details with respect to the target query. 2) **Diversity** indicates how rich and diverse the response is in providing various perspectives and insights related to the query. We employ GPT-40 as the judge to evaluate the model responses according to the grading criteria that gives scores ranging from 0 to 100 based on a two-step scoring scheme, as detailed in Appendix D.

For long narrative understanding, including NarrativeQA, Nocha, and Prelude, we uniformly use prediction accuracy (Acc) as the reported metric. Specifically, for NarrativeQA, prior studies (Bulian et al., 2022; Wang et al., 2024; Zhou et al., 2025) have shown that conventional token-level metrics such as Exact Match and F1 score usually fail to reflect actual semantic equivalence between hypothesis and reference answer, especially for abstractive answers. Therefore, we also apply GPT-40 for judging whether the LLM's prediction fully entails the reference answer, producing a binary True/False decision.

5 RESULTS AND ANALYSIS

5.1 OVERALL RESULTS

Table 1 reports the overall results across all evaluation tasks. Our HGMEM consistently outperforms both single-step and multi-step RAG baselines on every dataset. Importantly, our HGMEM with Qwen2.5-32B-Instruct matches or even outperforms baselines powered by the stronger GPT-4o, underscoring its value in resource-efficient scenarios.

The baselines exhibit mixed performance patterns reflecting their respective representational strengths. For instance, HippoRAG v2 relies on knowledge triples, which provides strong fact representation but limited coverage of events and plots. As a result, it performs well on NoCha but falls behind NaiveRAG on NarrativeQA. In contrast, GraphRAG and LightRAG excel at building global representations but are weaker at capturing fine-grained details, leading them to outperform other

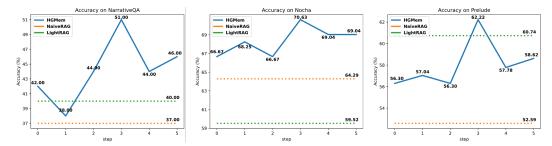


Figure 3: Prediction accuracies at different steps using Qwen2.5-32B-Instruct on long narrative understanding datasets.

Table 2: Ablation results using Qwen2.5-32B-Instruct. "w/. GE Only" and "w/. LI Only" stand for subquery generation strategies with only Global Exploration and Local Investigation, respectively. "w/o. Update" and "w/o. Merging" refer to HGMEM ablating update and merging operations during memory evolving, respectively.

Ablation Type	Method	Longbench		NarrativeQA	Nocha	Prelude	
ribiation Type	Memod	Comprehensiveness	Diversity	Acc (%)	Acc (%)	Acc (%)	
	НСМЕМ	64.18	66.51	51.00	70.63	62.22	
Subquery Strategy	w/. GE Only	59.25	61.67	47.00	68.25	59.26	
	w/. LI Only	61.38	63.82	43.00	63.49	60.00	
	НСМЕМ	64.18	66.51	51.00	70.63	62.22	
Memory Evolution	w/o. Update	62.48	64.92	50.00	68.25	60.00	
-	w/o. Merging	61.76	61.80	43.00	61.11	57.78	

baselines on Prelude and NarrativeQA. The two multi-step RAG methods, which mainly employ working memory to iteratively generate subqueries in a chaining fashion, struggle on sense-making questions, where integrating higher-order relationships is essential.

In comparison, our HGMEM provides strong compositional representations that span from facts to plots, equipping LLM reasoning with high-order correlations and integrated evidence. This enables it to meet the diverse requirements posed by the evaluation tasks.

5.2 Performance at Different Steps

During the execution of our multi-step RAG system, the memory progressively evolves and guides the LLM to proceed retrieval and reasoning. To inspect the effects of memory evolving over multiple interaction steps, we force the LLM to generate responses at every step for totally six turns, even if it originally decides to terminate iteration earlier. Figure 3 presents the performances at different steps using Qwen2.5-32B-Instruct on long narrative understanding tasks. Note that $t{=}0$ represents the initial step when the target query \hat{q} is used for retrieval. We can observe that our HGMEM achieves its best performance at $t{=}3$, mostly outperforming NaiveRAG and LightRAG baselines across steps. More steps bring no further improvements at higher cost.

5.3 ABLATION STUDIES

Subquery Generation Strategy. When the LLM determines to acquire more information from \mathcal{D} and \mathcal{G} , our HGMEM adopts an adaptive memory-based subquery generation strategy that generates subqueries either for focused local investigation and broad global exploration (Section 3.4). To investigate the effects of such strategy, in Table 2, we compare the adaptive strategy to the variants that involve only *Local Investigation* or *Global Exploration*, represented as "w/. LI Only" and "w/. GE Only" respectively. The results show that both "w/. LI Only" and "w/. GE Only" significantly underperform the adaptive strategy across all datasets, demonstrating the effectiveness and necessity of adaptively combining the two modes during subquery generation.

Effects of Update and Merging Operations. The memory evolving in our HGMEM involves update, insertion and merging operations, where merging is especially critical to building high-

Table 3: Average number of entities per hyperedge $(Avg-N_{\tilde{v}})$ in final memory and prediction accuracy (Acc) for a subset of 120 sampled primitive and sense-making queries.

Ouery Type	Method	NarrativeQA		Nocha		Prelude	
Query Type		Avg - $N_{\tilde{v}}$	Acc (%)	$\overline{\textbf{Avg-}N_{ ilde{v}}}$	Acc (%)	$\overline{\textbf{Avg-}N_{ ilde{v}}}$	Acc (%)
Primitive	HGMEM w/o. Merging	3.35 3.32	70.00 70.00	3.78 3.42	60.00 65.00	3.85 3.73	55.00 60.00
Sense-making	HGMEM w/o. Merging	7.07 4.10	40.00 30.00	7.97 3.80	70.00 60.00	5.25 3.74	60.00 55.00

order correlations from primitive facts. Because insertion is indispensable, we just carry out ablation experiments on all datasets using Qwen2.5-32B-Instruct to assess the effects of update and merging operations, as shown in Table 2. Compared to the "HGMEM", removing either operation leads to a performance drop, while removing merging ("w/o. Merging") causes a substantially larger degradation than removing update ("w/o. Update"). It reflects the effectiveness of both operations, especially highlighting the importance of high-order correlations built through merging operations.

5.4 Dissecting Query Resolving: Primitive vs. Sense-making

To better understand how our proposed HGMEM brings improvement to the evaluation tasks, we conduct a targeted analysis across different query types. Specifically, we randomly sample 40 queries from each long narrative understanding dataset used in our experiments, yielding a total of 120 queries. These are then manually categorized into two representative types:

- *Primitive Query*: Queries that primarily require locating directly associated chunks, which can often be resolved with local evidence and focus on straightforward factual information.
- Sense-making Query: Queries that require deeper comprehension by connecting and integrating multiple pieces of evidence, emphasizing the construction of higher-order relationships and interpretation beyond surface retrieval.

We compare both prediction accuracy and the average number of entities per hyperedge $(Avg-N_{\tilde{v}})$ in memory before generating final responses. The latter serves as a quantitative indicator of relationship complexity. Table 3 shows that on *sense-making queries*, our full "HGMEM" achieves higher accuracy with considerably larger $Avg-N_{\tilde{v}}$ than "HGMEM w/o. Merging", demonstrating that forming higher-order correlations enhances comprehension. In contrast, for primitive queries, "HGMEM" yields comparable or slightly lower accuracy relative to "HGMEM w/o. Merging". This is likely because the full model still tends to associate additional pieces of relevant evidence (as indicated by the slightly higher $Avg-N_{\tilde{v}}$), even though the primitive evidence alone is sufficient to answer straightforward queries, resulting in redundancy.

Notably, the $Avg-N_{\bar{v}}$ on sense-making queries consistently exceeds that on primitive queries, especially when merging is applied. Taken together, these results indicate that HGMEM improves contextual understanding by constructing high-order correlations for complex relational reasoning, rather than relying on shallow accumulation of surface facts.

6 Conclusion

In this work, we propose HGMEM, the hypergraph-based memory mechanism that aims at improving multi-step RAG by enabling the evolving of memory into more sophisticated forms for complex relational modeling. In HGMEM, the memory is structured as a hypergraph composed of a set of hyperedges as separate memory points. HGMEM allows the memory to progressively establish high-order correlations among previously accumulated primitive facts during the execution of multi-step RAG systems, guiding LLMs to organize and connect thoughts for a focal problem. Extensive experiments and in-depth analysis validate the effectiveness of our method over strong RAG baselines on challenging datasets featuring global sense-making questions over long context.

7 REPRODUCIBILITY STATEMENT

To ensure reproducibility, we introduce the usage and statistics of our used datasets in Section 4.1 and Appendix A. We also give the implementation details about the offline graph construction, system deployment and configuration in Section 4.2. Appendix C describes the procedures for subquery generation with prompts. Appendix D gives the evaluation prompts for scoring model responses in generative sense-making QA task.

REFERENCES

- Yushi Bai, Shangqing Tu, Jiajie Zhang, Hao Peng, Xiaozhi Wang, Xin Lv, Shulin Cao, Jiazheng Xu, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. Longbench v2: Towards deeper understanding and reasoning on realistic long-context multitasks. In *Proceedings of Association for Computational Linguistics*, pp. 3639–3664, 2025.
- Jannis Bulian, Christian Buck, Wojciech Gajewski, Benjamin Börschinger, and Tal Schuster. Tomayto, tomahto. beyond token-level answer equivalence for question answering evaluation. *CoRR*, abs/2202.07654, 2022.
- Howard Chen, Ramakanth Pasunuru, Jason Weston, and Asli Celikyilmaz. Walking down the memory maze: Beyond context limit through interactive reading. *CoRR*, abs/2310.05029, 2023.
- Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. BGE m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation. CoRR, abs/2402.03216, 2024.
- Prateek Chhikara, Dev Khant, Saket Aryan, Taranjeet Singh, and Deshraj Yadav. Mem0: Building production-ready AI agents with scalable long-term memory. *CoRR*, abs/2504.19413, 2025.
- Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. From local to global: A graph RAG approach to query-focused summarization. *CoRR*, abs/2404.16130, 2024.
- Masoomali Fatehkia, Ji Kim Lucas, and Sanjay Chawla. T-RAG: lessons from the LLM trenches. *CoRR*, abs/2402.07483, 2024.
- Yifan Feng, Haoxuan You, Zizhao Zhang, Rongrong Ji, and Yue Gao. Hypergraph neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 3558–3565, 2019.
- Xinyan Guan, Jiali Zeng, Fandong Meng, Chunlei Xin, Yaojie Lu, Hongyu Lin, Xianpei Han, Le Sun, and Jie Zhou. Deeprag: Thinking to retrieval step by step for large language models. *CoRR*, abs/2502.01142, 2025.
- Zirui Guo, Lianghao Xia, Yanhua Yu, Tu Ao, and Chao Huang. Lightrag: Simple and fast retrieval-augmented generation. *CoRR*, abs/2410.05779, 2024.
- Bernal Jimenez Gutierrez, Yiheng Shu, Yu Gu, Michihiro Yasunaga, and Yu Su. Hipporag: Neurobiologically inspired long-term memory for large language models. In *Proceedings of Neural Information Processing Systems*, 2024.
- Bernal Jiménez Gutiérrez, Yiheng Shu, Weijian Qi, Sizhe Zhou, and Yu Su. From RAG to memory: Non-parametric continual learning for large language models. *CoRR*, abs/2502.14802, 2025.
- Nicola Jones. Openai's' deep research'tool: is it useful for scientists? *Nature*, 2025.
- Marzena Karpinska, Katherine Thai, Kyle Lo, Tanya Goyal, and Mohit Iyyer. One thousand and one pairs: A "novel" challenge for long-context language models. In *Proceedings of EMNLP*, pp. 17048–17085, 2024.
- Tomás Kociský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis, and Edward Grefenstette. The narrativeqa reading comprehension challenge. *Transactions of the Association for Computational Linguistics*, 6:317–328, 2018.

- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the Symposium on Operating Systems Principles*, pp. 611–626, 2023.
 - Kuang-Huei Lee, Xinyun Chen, Hiroki Furuta, John F. Canny, and Ian Fischer. A human-inspired reading agent with gist memory of very long contexts. In *Proceedings of International Conference on Machine Learning*, 2024.
 - Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. CAMEL: communicative agents for "mind" exploration of large language model society. In *Proceedings of Neural Information Processing Systems*, 2023.
 - Shilong Li, Yancheng He, Hangyu Guo, Xingyuan Bu, Ge Bai, Jie Liu, Jiaheng Liu, Xingwei Qu, Yangguang Li, Wanli Ouyang, Wenbo Su, and Bo Zheng. Graphreader: Building graph-based agent to enhance long-context abilities of large language models. In *Findings of Empirical Methods in Natural Language Processing*, pp. 12758–12786, 2024.
 - Junru Lu, Siyu An, Mingbao Lin, Gabriele Pergola, Yulan He, Di Yin, Xing Sun, and Yunsheng Wu. Memochat: Tuning llms to use memos for consistent long-range open-domain conversation. *CoRR*, abs/2308.08239, 2023.
 - Barlas Oguz, Xilun Chen, Vladimir Karpukhin, Stan Peshterliev, Dmytro Okhonko, Michael Sejr Schlichtkrull, Sonal Gupta, Yashar Mehdad, and Scott Yih. Unik-qa: Unified representations of structured and unstructured knowledge for open-domain question answering. In *Findings of the Association for Computational Linguistics: NAACL*, pp. 1535–1546, 2022.
 - Hongjin Qian, Zheng Liu, Peitian Zhang, Kelong Mao, Defu Lian, Zhicheng Dou, and Tiejun Huang. Memorag: Boosting long context processing with global memory-enhanced retrieval augmentation. In *Proceedings of WWW 2025*, pp. 2366–2377, 2025.
 - Preston Rasmussen, Pavlo Paliychuk, Travis Beauvais, Jack Ryan, and Daniel Chalef. Zep: A temporal knowledge graph architecture for agent memory. *CoRR*, abs/2501.13956, 2025.
 - Parth Sarthi, Salman Abdullah, Aditi Tuli, Shubh Khanna, Anna Goldie, and Christopher D. Manning. RAPTOR: recursive abstractive processing for tree-organized retrieval. In *Proceedings of International Conference on Learning Representations*, 2024.
 - Minjie Shen and Qikai Yang. From mind to machine: The rise of manus AI as a fully autonomous digital agent. *CoRR*, abs/2505.02024, 2025.
 - Wenyu Tao, Xiaofen Xing, Yirong Chen, Linyi Huang, and Xiangmin Xu. Treerag: Unleashing the power of hierarchical storage for enhanced knowledge retrieval in long documents. In *Findings of the Association for Computational Linguistics*, pp. 356–371, 2025.
 - Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. In *Proceedings of the Association for Computational Linguistics*, pp. 10014–10037, 2023.
 - Juyuan Wang, Rongchen Zhao, Wei Wei, Yufeng Wang, Mo Yu, Jie Zhou, Jin Xu, and Liyan Xu. Comorag: A cognitive-inspired memory-organized RAG for stateful long narrative reasoning. *CoRR*, abs/2508.10419, 2025.
 - Yang Wang, Alberto Garcia Hernandez, Roman Kyslyi, and Nicholas Kersting. Evaluating quality of answers for retrieval-augmented generation: A strong LLM is all you need. *CoRR*, abs/2406.18064, 2024.
 - Wujiang Xu, Zujie Liang, Kai Mei, Hang Gao, Juntao Tan, and Yongfeng Zhang. A-MEM: agentic memory for LLM agents. *CoRR*, abs/2502.12110, 2025.
 - Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R. Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *Proceedings of International Conference on Learning Representations*, 2023.

Mo Yu, Tsz Ting Chung, Chulun Zhou, Tong Li, Rui Lu, Jiangnan Li, Liyan Xu, Haoshu Lu, Ning Zhang, Jing Li, and Jie Zhou. PRELUDE: A benchmark designed to require global comprehension and reasoning over long contexts. *CoRR*, abs/2508.09848, 2025.

- Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. Memorybank: Enhancing large language models with long-term memory. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 19724–19731, 2024.
- Chulun Zhou, Qiujing Wang, Mo Yu, Xiaoqian Yue, Rui Lu, Jiangnan Li, Yifan Zhou, Shunchi Zhang, Jie Zhou, and Wai Lam. The essence of contextual understanding in theory of mind: A study on question answering with story characters. In *Proceedings of the Association for Computational Linguistics*, pp. 22612–22631, 2025.

Table 4: Statistics of data used in our experiments. #Documents, Avg. #Tokens and #Queries represent the number of documents, average tokens per document and the total number of queries, respectively.

	Longbench (Financial)	Longbench (Governmental)	Longbench (Legal)	NarrativeQA	Nocha	Prelude
#Documents	20	22	7	10	4	5
Avg. #Tokens	266k	256k	194k	218k	139k	280k
#Queries	100	98	55	100	126	135

A DATASET STATISTICS

Generative Sense-making QA. We retain a portion of long documents with more than 100k tokens from Longbench V2 (Bai et al., 2025) which was originally comprised of six major task categories designed to assess the ability of LLMs to handle long-context problems. In our experiments, we select three domains of documents from the category of single-document QA, including Financial, Governmental and Legal.

Long Narrative Understanding. We use the following public benchmarks:

- NarrativeQA (Kociský et al., 2018): It is one of the most widely used benchmarks for story question answering. Because of its question construction strategy over high-level book summaries, the task places greater emphasis on synthesis and inference beyond local texts. In contrast, many other existing long-context QA tasks can often be solved with only local evidence, as shown by studies in (Yu et al., 2025). For evaluation, we randomly sample 10 long books exceeding 100k tokens, together with their associated queries, from the complete benchmark.
- NoCha (Karpinska et al., 2024): The task involves discriminating minimally different pairs of true and false claims about English fictional books. Although the format may appear different from sense-making questions, NoCha is explicitly designed to require constructing a global understanding of the book in relation to the focal statement. Since the official test set is hidden, we conduct experiments using only the publicly released subset.
- **Prelude** (Yu et al., 2025): This benchmark assesses LLMs' global comprehension and deep reasoning by requiring them to determine whether a character's prequel story is consistent with the original book. Most instances of this task demand integrating multiple pieces of evidence or even forming a holistic impression of the character's storyline. In our experiments, we use all English books included in Prelude for evaluation.

Table 4 gives the detailed statistics about the data used in our experiments, including the number of documents, average tokens per document and the total number of queries. Generative sense-making QA task involves documents from Longbench V2 benchmark in *Financial*, *Government* and *Legal* domains. Long narrative understanding task uses NarrativeQA, Nocha and Prelude benchmarks.

B COMPARISON BASELINES

In our experiments, we compare our methods to traditional RAG and Multi-step RAG methods. Traditional RAG includes:

- NaiveRAG just uses the target query to retrieve a set of text chunks from the document for dealing with queries.
- **GraphRAG** (Edge et al., 2024) constructs knowledge graph from plain-text documents and build a hierarchy of communities of closely related entities before using an LLM to make responses.
- LightRAG (Guo et al., 2024) also builds a graph structure and employs a dual-level retrieval strategy from both low-level and high-level evidence discovery.
- HippoRAG v2 (Gutiérrez et al., 2025) creates a knowledge graph and adopts the Personalized PageRank algorithm with dense-sparse integration of passages into the graph search process for resolving queries.

```
You are an intelligent assistant responsible for dealing with the [Main Query] by making appropriate operations as specified.

With respect to the [Main Query], you have consolidated some memory points in your [Memory] describing what you have already known regarding the [Main Query].

Each memory point can be seen as a specific aspect relevant to the [Main Query], providing necessary details or insights from its perspective.
702
703
704
                                                                                         Your task is to analyze the [Main Query] and [Memory], then determine whether current [Memory] has been sufficient to comprehensively resolve the [Main Query] If not sufficient, you need to indicate what you want to further investigate.
705
                                                                                          Step 1.

Make appropriate judgement following the logic branches below.

Case 1: If the [Memory] has been sufficient to completely resolve the [Main Query], output <None> in [Concerns].

Case 1: If the [Memory] is not sufficient, determine current situation should be attributed to which of the following subcases.

Case 2: If the [Memory] is not sufficient, determine current situation should be attributed to which of the following subcases.

Case 2: There are some specific memory points which you want to further investigate more details about.

Case 2: There are unexplored aspects that go beyond the scope of current [Memory] (i.e. not related to any of the existing memory points).
706
708
                                                                                           Step 2. Output as **Example of Anticipated Output Format**.

Specifically, give your judgement in Judgement] using corresponding case index (1, 2.1 or 2.2).

Then, generate several concerns that aim at exploring details or aspects not addressed by current [Memory] to better resolve the [Main Query]

When case 2.1, generate up to [num_concerns] concerns, each of which targets at a specific memory point. For each concern, specify the index of its corresponding memory
709
710
                                                                                          point.

When case 2.2, generate up to {num_concerns} concerns that probe meaningful information not yet covered by current [Memory]
711
                                                                                                                           mple of Anticipated Output Format for Case 1-##########
712
713
                                                                                          714
715
                                                                                           2{tuple_delimiter}your_concern_3{record_delimiter}
{completion_delimiter}
716
                                                                                           ###########Example of Anticipated Output Format for Case 2.2-############
                                                                                         717
718
719
720
                                                                                          721
                                                                                         (nemony):
(memony):
(memony):
"Mote that:
"Note that:
(1) Your concern should be concise and suggest what further details or aspect you subsequently will seek for.
(2) Only output the judgement, concerns, and the indices of corresponding memory points without any other of
(3) If current [Memory ] has covered most relevant perspectives, generate fewer concerns to avoid redundancy.
(4) Your generated concerns should be separated by "fecond_delimiter)".
722
723
724
725
                                                                                          726
```

Figure 4: The prompt for raising concerns either targeting at specific memory points or probing useful information outside current memory.

Multi-step RAG includes:

- **DeepRAG** (Guan et al., 2025) conducts multi-step reasoning as a Markov Decision Process by iteratively decomposing queries.
- ComoRAG (Wang et al., 2025) undergoes multi-step interactions with external data sources with a dynamic memory workspace, iteratively generateing probing queries and integrating the retrieved evidence into a global memory pool.

C SUBQUERY GENERATION

In Section 3.4, we introduce our designed adaptive memory-based subquery generation strategy. Specifically, the LLM analyzes the memory and determines to generate subqueries for local investigation or global exploration for better dealing with the target query. First, it raises relevant concerns that either target at specific memory points or aim at probing useful information outside current memory. Then, the LLM generates corresponding subqueries according to different cases. The prompts for raising concerns and generating subqueries are given in Figure 4 and Figure 5, respectively.

774 775 776

777 778

779

780

781

782

783 784

785 786

787

789

791

793

794

796

797

798

799

800

801 802

804

809

```
756
                                   You are an assistant responsible for dealing with the [Main Query].

Although you have had some relevant information in your [Memory], your current [Memory] is still not sufficient to comprehensively
                                    resolve the [Main Query] due to the concern given in [Concern].
                                    Therefore, you need to generate a subquery that aims at either retrieving more evidences or investigating unexplored aspects in
758
                                   [Subquery] to better deal with the [Main Query] ultimately.
                                    [Previous Subqueries] records a series of previous subqueries that have been raised before
760
                                    #########-Anticipated Output Format-#########
761
762
                                    [Main Query]: {query}
                                   [Memory]
764
765
                                    [Concern]
766
                                    {concern}
                                    [Previous Subqueries]
                                    {history_subqueries}
768
                                   *******************
769
                                     Note that:
                                    (1) Your generated subquery should be concise and address the concerns in your [Concern].
770
                                    (2) You should avoid generating a subquery that is overly similar to any one of the [Previous Subqueries] or [Main Query].
                                    (3) Only output your subquery without any other redundant content such as markup strings.
772
                                   Output:
```

Figure 5: The prompt for generating subqueries based on previously raised concerns.

D EVALUATION PROMPTS FOR GENERATIVE SENSE-MAKING QA

Given a [Paragraph] and a [Question], you will evaluate the quality of the [Response] in terms of Comprehensiveness

For the evaluation of generative sense-making QA, we leverage GPT-40 as an evaluator to assess the quality of model responses. Given the target query and the source paragraph from which the query was originated, the GPT-40 evaluator first indicates the level of comprehensiveness/diversity and then gives a final score within the value range of corresponding level. Detailed prompts for such LLM-as-a-Judge evaluation. Figure 6 and Figure 7 give the prompts for scoring the comprehensiveness and diversity, respectively.

```
[Paragraph]:{paragraph}
[Question]: {question}
[Response]:{response}
Comprehensiveness measures whether the [Response] comprehensively covers all key aspects in the [Paragraph] with respect to the
Level | score range | description
Level 1 | 0-20 | The response is extremely one-sided, leaving out key parts or important aspects of the question.
Level 2 | 20-40 | The response has some content, but it misses many important aspects of the question and is not comprehensive enough. Level 3 | 40-60 | The response is moderately comprehensive, covering the main aspects of the question, but there are still some omissions
Level 4 | 60-80 | The response is comprehensive, covering most aspects of the question, with few omissions.
Level 5 | 80-100 | The response is extremely comprehensive, covering almost all aspects of the question no omissions, enabling the reader to
gain a complete and thorough understanding.

Evaluate the [Response] using the criteria listed above, give a level of comprehensiveness in [Level] based on the description of the indicator,
then give a score in [Score] based on the corresponding value range, and finally explain in [Explanation]
(1) You should reference to the [Paragraph] and avoid misinterpreting any content of [Paragraph] as part of the [Response].
(2) Avoid excessively concerning very specific details. When the response mentions an aspect without providing very specific details, you
should consider this aspect as validly covered, as long as the omitted detail is not crucial to particularly mention with respect to the
[Question] in the whole scope of the response.
(3) If [Response] contains extra content not directly included in the [Paragraph], as long as the extra content is correct, do not consider the
extra content as defects for giving final evaluation.
(4) You should conform to the -Anticipated Output Format- and give your evaluation results in [Your Evaluation].
      [Level]: A level ranging from 1 to 5 # This should be a single number, not a range.

[Score]: A value ranging from 0 to 100 # This should be a single number satisfying the ranging constraint of the corresponding [Level], not a
range.
[Explanation]: xxx
[Your Evaluation]
```

Figure 6: The prompt for evaluating the comprehensiveness of a model response.

Given a [Paragraph] and a [Question], you will evaluate the quality of the [Response] in terms of Diversity. [Paragraph]: {paragraph} [Question]: {question} [Response]: {response} Diversity measures how varied and rich is the response in offering different perspectives and insights related to the question Level | score range | description | Level 1 | 0-20 | The response is extremely narrow and repetitive, providing only a single perspective or insight without exploring alternative viewpoints or additional information. Level 2 | 20-40 | The response offers a few different perspectives but remains largely superficial. It may touch on alternative viewpoints but does not elaborate or provide substantial insights.

Level 3 | 40-60 | The response moderately presents several perspectives with moderate depth. It begins to integrate different viewpoints and insights but may still miss some important angles or lack thorough exploration.

Level 4 | 60-80 | The response is rich in perspectives and insights. It basically explores multiple viewpoints and provides substantial evidence and examples to support each angle.

Level 5 | 80-100 | The response is exceptionally varied and rich in perspectives and insights. It offers a comprehensive exploration of the question, addressing multiple angles with depth and originality.

Evaluate the [Response] using the criteria listed above, give a level of comprehensiveness in [Level] based on the description of the indicator, then give a score in [Score] based on the corresponding value range, and finally explain in [Explanation]. (1) You should reference to the [Paragraph] and avoid misinterpreting any content of [Paragraph] as part of the [Response].

(2) If [Response] contains extra content not directly included in the [Paragraph], as long as the extra content is correct, do not consider the extra content as defects for giving final evaluation. (3) You should conform to the -Anticipated Output Format- and give your evaluation results in [Your Evaluation] [Score]: A value ranging from 0 to 100 # This should be a single number satisfying the ranging constraint of the corresponding [Level], not a [Explanation]: xxx [Your Evaluation]:

Figure 7: The prompt for evaluating the comprehensiveness of a model response.