LightRAG: Simple and Fast Retrieval-Augmented Generation

Zirui Guo¹, Lianghao Xia¹, Yanhua Yu², Tu Ao², Chao Huang^{1*}

University of Hong Kong¹
Beijing University of Posts and Telecommunications²
zrguo101@hku.hk aka_xia@foxmail.com chaohuang75@gmail.com

Abstract

Retrieval-Augmented Generation (RAG) systems enhance large language models (LLMs) by integrating external knowledge sources, enabling more accurate and contextually relevant responses tailored to user needs. However, existing RAG systems have significant limitations, including reliance on flat data representations and inadequate contextual awareness, which can lead to fragmented answers that fail to capture complex interdependencies. To address these challenges, we propose LightRAG, a novel framework that incorporates graph structures into text indexing and retrieval processes. This innovative approach employs a dual-level retrieval system that enhances comprehensive information retrieval from both lowand high-level knowledge discovery. Additionally, the integration of graph structures with vector representations facilitates efficient retrieval of related entities and their relationships, significantly improving response times while maintaining contextual relevance. This capability is further enhanced by an incremental update algorithm that ensures the timely integration of new data, allowing the system to remain effective and responsive in rapidly changing data environments. Extensive experimental validation demonstrates considerable improvements in retrieval accuracy and efficiency compared to existing approaches. LightRAG is publicly available as an open-source framework at: https://github.com/HKUDS/LightRAG.

1 Introduction

Retrieval-Augmented Generation (RAG) systems have been developed to enhance large language models (LLMs) by integrating external knowledge sources (Sudhi et al., 2024; Es et al., 2024; Salemi et al., 2024). This innovative integration allows LLMs to generate more accurate and contextually relevant responses, significantly improving their

*Corresponding Author: Chao Huang

utility in real-world applications. By adapting to specific domain knowledge (Tu et al., 2024), RAG systems ensure that the information provided is not only pertinent but also tailored to the user's needs. Furthermore, they offer access to up-to-date information (Zhao et al., 2024), which is crucial in rapidly evolving fields. Chunking plays a vital role in facilitating the retrieval-augmented generation process (Lyu et al., 2024). By breaking down a large external text corpus into smaller, more manageable segments, chunking significantly enhances the accuracy of information retrieval. This enables more targeted similarity searches, ensuring that the retrieved content is directly relevant to user queries.

However, existing RAG systems have key limitations that hinder their performance. First, many methods rely on flat data representations, restricting their ability to understand and retrieve information based on intricate relationships between entities. **Second**, these systems often lack the contextual awareness needed to maintain coherence across various entities and their interrelations, resulting in responses that may not fully address user queries. For example, consider a user asking, "How does the rise of electric vehicles influence urban air quality and transportation infrastructure?" Existing RAG methods might retrieve separate documents on electric vehicles, air pollution, and transportation challenges but struggle to synthesize them into a cohesive response. They may fail to explain how the adoption of electric vehicles can improve air quality, which in turn could affect public transportation planning. As a result, the answer may be fragmented and does not adequately capture the complex inter-dependencies among these topics.

To address these limitations, we propose incorporating graph structures into text indexing and relevant information retrieval. Graphs are particularly effective at representing the interdependencies among different entities (Rampášek et al., 2022), which enables a more nuanced understanding of re-

lationships. The integration of graph-based knowledge structures facilitates the synthesis of information from multiple sources into coherent and contextually rich responses. Despite these advantages, developing a fast and scalable graph-empowered RAG system that efficiently handles varying query volumes is crucial. In this work, we achieve an effective and efficient RAG system by addressing three key challenges: i) Comprehensive Information Retrieval. Ensuring comprehensive information retrieval that captures the full context of inter-dependent entities from all documents; ii) Enhanced Retrieval Efficiency. Improving retrieval efficiency over the graph-based knowledge structures to significantly reduce response times; iii) Rapid Adaptation to New Data. Enabling quick adaptation to new data updates, ensuring the system remains relevant in dynamic environments.

In response to the outlined challenges, we propose LightRAG, a model that seamlessly integrates a graph-based text indexing paradigm with a dual-level retrieval framework. This innovative approach enhances the system's capacity to capture complex inter-dependencies among entities, resulting in more coherent and contextually rich responses. LightRAG employs efficient dual-level retrieval strategies: low-level retrieval, which focuses on precise information about specific entities and their relationships, and high-level retrieval, which encompasses broader topics and themes. By combining both detailed and conceptual retrieval, LightRAG effectively accommodates a diverse range of quries, ensuring that users receive relevant and comprehensive responses tailored to their specific needs. Additionally, by integrating graph structures with vector representations, our framework facilitates efficient retrieval of related entities and relations while enhancing the comprehensiveness of results through relevant structural information from the constructed knowledge graph.

In summary, this work's key contributions are:

- General Aspect. We emphasize the importance
 of developing a graph-empowered RAG system
 to overcome the limitations of existing methods.
 By integrating graph structures into text indexing, we can effectively represent complex interdependencies among entities, fostering a nuanced understanding of relationships and enabling coherent, contextually rich responses.
- Methodologies. To enable an efficient and adaptive RAG system, we propose LightRAG, which

integrates a dual-level retrieval paradigm with graph-enhanced text indexing. This approach captures both low-level and high-level information for comprehensive, cost-effective retrieval. Without the need to rebuild the entire index, LightRAG reduces computational costs and accelerates adaptation, while its incremental update algorithm ensures timely integration of new data, maintaining efficacy in dynamic environments.

• Experimental Findings. Extensive experiments were conducted to evaluate the effectiveness of LightRAG in comparison to existing RAG models. These assessments focused on several key dimensions, including retrieval accuracy, model ablation, response efficiency, and adaptability to new information. The results demonstrated significant improvements over baseline methods.

2 Retrieval-Augmented Generation

Retrieval-Augmented Generation (RAG) integrates user queries with a collection of pertinent documents sourced from an external knowledge database, incorporating two essential elements: the **Retrieval Component** and the **Generation Component**. 1) The retrieval component is responsible for fetching relevant documents or information from the external knowledge database. It identifies and retrieves the most pertinent data based on the input query. 2) After the retrieval process, the generation component takes the retrieved information and generates coherent, contextually relevant responses. It leverages powerful LLMs to produce meaningful outputs. Formally, this RAG framework, denoted as \mathcal{M} , can be defined as follows:

$$\mathcal{M} = \left(\mathcal{G}, \ \mathcal{R} = (\varphi, \psi)\right),$$

$$\mathcal{M}(q; \mathcal{D}) = \mathcal{G}\left(q, \psi(q; \hat{\mathcal{D}})\right), \ \hat{\mathcal{D}} = \varphi(\mathcal{D}) \quad (1)$$

In this framework, $\mathcal G$ and $\mathcal R$ represent the generation module and the retrieval module, respectively, while q denotes the input query and D refers to the external database. The retrieval module $\mathcal R$ includes two key functionalities: i) **Data Indexer** $\varphi(\cdot)$: which involves building a specific data structure $\hat{\mathcal D}$ based on the external database D. ii) **Data Retriever** $\psi(\cdot)$: The relevant documents are obtained by comparing the query against the indexed data, also denoted as "relevant documents". By leveraging the information retrieved through $\psi(\cdot)$ along with the initial query q, the generative model

 $\mathcal{G}(\cdot)$ efficiently produces high-quality responses. This work targets several key points essential for an efficient and effective RAG system as follows:

- Comprehensive Information Retrieval: The indexing function $\varphi(\cdot)$ must be adept at extracting global information, as this is crucial for effective query answering using LLMs.
- Efficient and Low-Cost Retrieval: The indexed data structure $\hat{\mathcal{D}}$ must enable rapid and cost-efficient information retrieval to effectively handle a high volume of user queries.
- Fast Adaptation to Data Changes: The ability to swiftly and efficiently adjust the data structure to incorporate new information from the external knowledge base, is crucial for ensuring that the system remains current and relevant in an everchanging information landscape.

3 The LightRAG Architecture

3.1 Graph-based Text Indexing

Graph-Enhanced Entity and Relationship Extraction. Our LightRAG enhances the retrieval system by segmenting documents into smaller, more manageable pieces. This strategy allows for quick identification and access to relevant information without analyzing entire documents. Next, we leverage LLMs to identify and extract various entities (e.g., names, dates, locations, and events) along with the relationships between them. The information collected through this process will be used to create a comprehensive knowledge graph that highlights the connections and insights across the entire collection of documents. We formally represent this graph generation module as follows:

$$\hat{\mathcal{D}} = (\hat{\mathcal{V}}, \hat{\mathcal{E}}) = \text{Dedupe} \circ \text{Prof}(\mathcal{V}, \mathcal{E}),$$

$$\mathcal{V}, \mathcal{E} = \cup_{\mathcal{D}_i \in \mathcal{D}} \text{Recog}(\mathcal{D}_i)$$
 (2)

where $\hat{\mathcal{D}}$ represents the resulting knowledge graphs. To generate this data, we apply three main processing steps to the raw text documents \mathcal{D}_i . These steps utilize a LLM for text analysis and processing. Details about the prompt templates and specific settings for this part can be found in Appendix 9.4.2. The functions used in our graph-based text indexing paradigm are described as:

• Extracting Entities and Relationships. $R(\cdot)$: This function prompts a LLM to identify entities

(nodes) and their relationships (edges) within the text data. For instance, it can extract entities like "Cardiologists" and "Heart Disease," and relationships such as "Cardiologists diagnose Heart Disease" from the text: "Cardiologists assess symptoms to identify potential heart issues." To improve efficiency, the raw text \mathcal{D} is segmented into multiple chunks \mathcal{D}_i .

• LLM Profiling for Key-Value Pair Generation. P(·): We employ a LLM-empowered profiling function, P(·), to generate a text key-value pair (K, V) for each entity node in V and relation edge in E. Each index key is a word or short phrase that enables efficient retrieval, while the corresponding value is a text paragraph summarizing relevant snippets from external data to aid in text generation. Entities use their names as the sole index key, whereas relations may have mul-

tiple keys derived from LLM enhancements that

include global themes from connected entities.

 Deduplication to Optimize Graph Operations. D(·): Finally, we implement a deduplication function, D(·), that identifies and merges identical entities and relations from different segments of the raw text D_i. This process effectively reduces the overhead associated with graph operations on D̂ by minimizing the graph's size, leading to more efficient data processing.

Our LightRAG offers two advantages through its graph-based text indexing paradigm. *First*, **Comprehensive Information Understanding**. The constructed graph enables the extraction of global information from multi-hop subgraphs, greatly enhancing LightRAG's ability to handle complex queries that span multiple document chunks. *Second*, **Enhanced Retrieval Performance**. the key-value data structures derived from the graph are optimized for rapid and precise retrieval. This provides a superior alternative to less accurate embedding matching methods (Gao et al., 2023) and inefficient chunk traversal techniques (Edge et al., 2024) commonly used in existing approaches.

Fast Adaptation to Incremental Knowledge Base. To efficiently adapt to evolving data changes while ensuring accurate and relevant responses, our LightRAG incrementally updates the knowledge base without the need for complete reprocessing of the entire external database. For a new document \mathcal{D}' , the incremental update algorithm processes it using the same graph-based indexing steps φ as

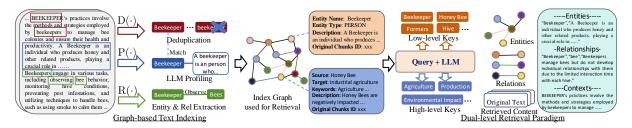


Figure 1: Overall architecture of the proposed LightRAG framework.

before, resulting in $\hat{\mathcal{D}}' = (\hat{\mathcal{V}}', \hat{\mathcal{E}}')$. Subsequently, LightRAGcombines the new graph data with the original by taking the union of the node sets $\hat{\mathcal{V}}$ and $\hat{\mathcal{V}}'$, as well as the edge sets $\hat{\mathcal{E}}$ and $\hat{\mathcal{E}}'$.

Two key objectives guide our approach to fast adaptation for the incremental knowledge base: **Seamless Integration of New Data**. By applying a consistent methodology to new information, the incremental update module allows the LightRAG to integrate new external databases without disrupting the existing graph structure. This approach preserves the integrity of established connections, ensuring that historical data remains accessible while enriching the graph without conflicts or redundancies. Reducing Computational Over**head** . By eliminating the need to rebuild the entire index graph, this method reduces computational overhead and facilitates the rapid assimilation of new data. Consequently, LightRAG maintains system accuracy, provides current information, and conserves resources, ensuring users receive timely updates and enhancing the RAG effectiveness.

3.2 Dual-level Retrieval Paradigm

To retrieve relevant information from both specific document chunks and their complex interdependencies, our LightRAG proposes generating query keys at both detailed and abstract levels.

- **Specific Queries**. These queries are detailoriented and typically reference specific entities within the graph, requiring precise retrieval of information associated with particular nodes or edges. For example, a specific query might be, "Who wrote 'Pride and Prejudice'?"
- **Abstract Queries**. In contrast, abstract queries are more conceptual, encompassing broader topics, summaries, or overarching themes that are not directly tied to specific entities. An example of an abstract query is, "How does artificial intelligence influence modern education?"

To accommodate diverse query types, the LightRAG employs two distinct retrieval strategies

within the dual-level retrieval paradigm. This ensures that both specific and abstract inquiries are addressed effectively, allowing the system to deliver relevant responses tailored to user needs.

- Low-Level Retrieval. This level is primarily focused on retrieving specific entities along with their associated attributes or relationships. Queries at this level are detail-oriented and aim to extract precise information about particular nodes or edges within the graph.
- High-Level Retrieval. This level addresses broader topics and overarching themes. Queries at this level aggregate information across multiple related entities and relationships, providing insights into higher-level concepts and summaries rather than specific details.

Integrating Graph and Vectors for Efficient Retrieval. By combining graph structures with vector representations, the model gains a deeper insight into the interrelationships among entities. This synergy enables the retrieval algorithm to effectively utilize both local and global keywords, streamlining the search process and improving the relevance of results.

- (i) **Query Keyword Extraction**. For a given query q, the retrieval algorithm of LightRAG begins by extracting both local query keywords $k^{(l)}$ and global query keywords $k^{(g)}$.
- (ii) Keyword Matching. The algorithm uses an efficient vector database to match local query keywords with candidate entities and global query keywords with relations linked to global keys.
- (iii) Incorporating High-Order Relatedness. To enhance the query with higher-order relatedness, LightRAGfurther gathers neighboring nodes within the local subgraphs of the retrieved graph elements. This process involves the set $\{v_i|v_i\in\mathcal{V}\land(v_i\in\mathcal{N}_v\vee v_i\in\mathcal{N}_e)\}$, where \mathcal{N}_v and \mathcal{N}_e represent the one-hop neighboring nodes of the retrieved nodes v and edges e, respectively.

This dual-level retrieval paradigm of LightRAG not only facilitates efficient retrieval of related entities and relations through keyword matching, but also enhances the comprehensiveness of retrieval results by integrating relevant structural information from the constructed knowledge graph.

3.3 Retrieval-Augmented Answer Generation

Utilization of Retrieved Information. Utilizing the retrieved information $\psi(q; \hat{\mathcal{D}})$, our LightRAG employs a general-purpose LLM to generate answers based on the collected data. This data comprises concatenated values V from relevant entities and relations, produced by the profiling function $P(\cdot)$. It includes names, descriptions of entities and relations, and excerpts from the original text.

Context Integration and Answer Generation. By unifying the query with this multi-source text, the LLM generates informative answers tailored to the user's needs, ensuring alignment with the query's intent. This approach streamlines the answer generation process by integrating both context and query into the LLM model, as illustrated in detailed examples (Appendix 9.3).

3.4 Complexity Analysis of LightRAG

In this section, we analyze the complexity of our proposed LightRAG framework, which can be divided into two main parts. The first part is the graph-based Index phase. During this phase, we use the large language model (LLM) to extract entities and relationships from each chunk of text. As a result, the LLM needs to be called total tokens chunk size times. Importantly, there is no additional overhead involved in this process, making our approach highly efficient in managing updates to new text.

The second part of the process involves the graph-based retrieval phase. For each query, we first utilize the large language model (LLM) to generate relevant keywords. Similar to current Retrieval-Augmented Generation (RAG) systems (Gao et al., 2023, 2022; Chan et al., 2024), our retrieval mechanism relies on vector-based search. However, instead of retrieving chunks as in conventional RAG, we concentrate on retrieving entities and relationships. This approach markedly reduces retrieval overhead compared to the community-based traversal method used in GraphRAG.

4 Evaluation

We conduct empirical evaluations on benchmark data to assess the effectiveness of the proposed LightRAG framework by addressing the following research questions: • (RQ1): How does LightRAG compare to existing RAG baseline methods in terms of generation performance? • (RQ2): How do dual-level retrieval and graph-based indexing enhance the generation quality of LightRAG? • (RQ3): What are the costs associated with LightRAG, as well as its adaptability to data changes?

4.1 Experimental Settings

To evaluate the effectiveness of LightRAG, we conducted experiments on four datasets from the Ultra-Domain benchmark (Qian et al., 2024), covering diverse domains such as Agriculture, CS, Legal, and Mixed. Each dataset ranges from 600,000 to 5,000,000 tokens. We compare LightRAG with state-of-the-art RAG methods, including Naive RAG, RQ-RAG, HyDE, and GraphRAG. Using a robust LLM-based evaluation, we assess performance across four dimensions: Comprehensiveness, Diversity, Empowerment, and Overall.

Detailed descriptions of datasets, question generation, and baselines are provided in Appendix 9.1.

4.2 RAG Performance Comparison (RQ1)

We compare LightRAG against each baseline across various evaluation dimensions and datasets. The results are presented in Table 1. Based on these findings, we draw the following conclusions:

Graph-Enhanced RAG Superiority in Large-Scale Corpora. When handling large token counts and complex queries requiring comprehensive dataset understanding, graph-based RAG systems like LightRAG and GraphRAG consistently outperform chunk-based methods such as NaiveRAG, HyDE, and RQRAG. This performance gap widens with dataset size—in the largest dataset (Legal), baseline methods achieve only 20% win rates compared to LightRAG's dominance. This trend highlights the advantages of graph-enhanced RAG systems in capturing complex semantic dependencies within large-scale corpora, enabling better knowledge understanding and improved generalization.

Enhancing Response Diversity with LightRAG. Compared to various baselines, LightRAG demonstrates a significant advantage in the Diversity metric, particularly within the larger Legal dataset. Its consistent lead in this area underscores LightRAG's effectiveness in generating a wider range of responses, especially in scenarios where diverse content is essential. We attribute

Table 1: Win rates (%) of baselines v.s. LightRAG across four datasets and four evaluation dimensions.

	Agriculture		C	2S	Legal		Mix	
	NaiveRAG	LightRAG	NaiveRAG	LightRAG	NaiveRAG	LightRAG	NaiveRAG	LightRAG
Comprehensiveness	32.4%	<u>67.6%</u>	38.4%	61.6%	16.4%	83.6%	38.8%	61.2%
Diversity	23.6%	<u>76.4%</u>	38.0%	62.0%	13.6%	86.4%	32.4%	<u>67.6%</u>
Empowerment	32.4%	<u>67.6%</u>	38.8%	61.2%	16.4%	83.6%	42.8%	<u>57.2%</u>
Overall	32.4%	<u>67.6%</u>	38.8%	<u>61.2%</u>	15.2%	84.8%	40.0%	60.0%
	RQ-RAG	LightRAG	RQ-RAG	LightRAG	RQ-RAG	LightRAG	RQ-RAG	LightRAG
Comprehensiveness	31.6%	68.4%	38.8%	61.2%	15.2%	84.8%	39.2%	60.8%
Diversity	29.2%	<u>70.8%</u>	39.2%	60.8%	11.6%	88.4%	30.8%	<u>69.2%</u>
Empowerment	31.6%	<u>68.4%</u>	36.4%	<u>63.6%</u>	15.2%	84.8%	42.4%	<u>57.6%</u>
Overall	32.4%	<u>67.6%</u>	38.0%	<u>62.0%</u>	14.4%	<u>85.6%</u>	40.0%	60.0%
	HyDE	LightRAG	HyDE	LightRAG	HyDE	LightRAG	HyDE	LightRAG
Comprehensiveness	26.0%	74.0%	41.6%	58.4%	26.8%	73.2%	40.4%	59.6%
Diversity	24.0%	76.0%	38.8%	61.2%	20.0%	80.0%	32.4%	<u>67.6%</u>
Empowerment	25.2%	74.8%	40.8%	<u>59.2%</u>	26.0%	74.0%	46.0%	54.0%
Overall	24.8%	<u>75.2%</u>	41.6%	<u>58.4%</u>	26.4%	73.6%	42.4%	<u>57.6%</u>
	GraphRAG	LightRAG	GraphRAG	LightRAG	GraphRAG	LightRAG	GraphRAG	LightRAG
Comprehensiveness	45.6%	54.4%	48.4%	51.6%	48.4%	51.6%	50.4%	49.6%
Diversity	22.8%	<u>77.2%</u>	40.8%	<u>59.2%</u>	26.4%	<u>73.6%</u>	36.0%	64.0%
Empowerment	41.2%	<u>58.8%</u>	45.2%	<u>54.8%</u>	43.6%	<u>56.4%</u>	50.8%	49.2%
Overall	45.2%	<u>54.8%</u>	48.0%	<u>52.0%</u>	47.2%	<u>52.8%</u>	<u>50.4%</u>	49.6%

this advantage to LightRAG's dual-level retrieval paradigm, which facilitates comprehensive information retrieval from both low-level and high-level dimensions. This approach effectively leverages graph-based text indexing to consistently capture the full context in response to queries.

LightRAG's Superiority over GraphRAG:

While both LightRAG and GraphRAG use graph-based retrieval mechanisms, LightRAG consistently outperforms GraphRAG, particularly in larger datasets with complex language In the Agriculture, CS, and Lecontexts. gal datasets—each containing millions of tokens-LightRAG shows a clear advantage, significantly surpassing GraphRAG and highlighting its strength in comprehensive information understanding within diverse environments. Enhanced Response Variety: By integrating low-level retrieval of specific entities with high-level retrieval of broader topics, LightRAG boosts response diversity. This dual-level mechanism effectively addresses both detailed and abstract queries, ensuring a thorough grasp of information. Complex Query **Handling**: This approach is especially valuable in scenarios requiring diverse perspectives. By accessing both specific details and overarching themes, LightRAG adeptly responds to complex queries involving interconnected topics, providing contextually relevant answers of high quality.

4.3 Ablation Studies (RQ2)

We also conduct ablation studies to evaluate the impact of our dual-level retrieval paradigm and the effectiveness of our graph-based text indexing in LightRAG. The results are presented in Table 2.

Effectiveness of Dual-level Retrieval Paradigm. We begin by analyzing the effects of low-level and high-level retrieval paradigms. We compare two ablated models against LightRAG across four datasets. Here are our key observations:

- Low-level-only Retrieval: The -High variant removes high-order retrieval, leading to a significant performance decline across nearly all datasets and metrics. This drop is mainly due to its emphasis on the specific information, which focuses excessively on entities and their immediate neighbors. While this approach enables deeper exploration of directly related entities, it struggles to gather information for complex queries that demand comprehensive insights.
- High-level-only Retrieval: The -Low variant prioritizes capturing a broader range of content by leveraging entity-wise relationships rather than focusing on specific entities. This approach offers a significant advantage in comprehensiveness, allowing it to gather more extensive and varied information. However, the trade-off is a reduced depth in examining specific entities, which can limit its ability to provide highly de-

				~ .
Table 2: Performance	of ablated version	s of LightRA(+	using Naivek	A(+ as reference
radic 2. I cirolinance	or abrated version	o oi Ligitua io.	, using ranver	as reference.

	Agriculture		(CS	Le	gal	Mix	
	NaiveRAG	LightRAG	NaiveRAG	LightRAG	NaiveRAG	LightRAG	NaiveRAG	LightRAG
Comprehensiveness	32.4%	67.6%	38.4%	61.6%	16.4%	83.6%	38.8%	61.2%
Diversity	23.6%	76.4%	38.0%	62.0%	13.6%	86.4%	32.4%	<u>67.6%</u>
Empowerment	32.4%	67.6%	38.8%	61.2%	16.4%	83.6%	42.8%	<u>57.2%</u>
Overall	32.4%	<u>67.6%</u>	38.8%	61.2%	15.2%	84.8%	40.0%	60.0%
	NaiveRAG	-High	NaiveRAG	-High	NaiveRAG	-High	NaiveRAG	-High
Comprehensiveness	34.8%	65.2%	42.8%	57.2%	23.6%	<u>76.4%</u>	40.4%	<u>59.6%</u>
Diversity	27.2%	72.8%	36.8%	63.2%	16.8%	83.2%	36.0%	64.0%
Empowerment	36.0%	64.0%	42.4%	<u>57.6%</u>	22.8%	<u>77.2%</u>	47.6%	<u>52.4%</u>
Overall	35.2%	64.8%	44.0%	56.0%	22.0%	<u>78.0%</u>	42.4%	57.6%
	NaiveRAG	-Low	NaiveRAG	-Low	NaiveRAG	-Low	NaiveRAG	-Low
Comprehensiveness	36.0%	64.0%	43.2%	56.8%	19.2%	80.8%	36.0%	64.0%
Diversity	28.0%	72.0%	39.6%	60.4%	13.6%	86.4%	33.2%	66.8%
Empowerment	34.8%	<u>65.2%</u>	42.8%	<u>57.2%</u>	16.4%	83.6%	35.2%	<u>64.8%</u>
Overall	34.8%	65.2%	43.6%	<u>56.4%</u>	18.8%	81.2%	35.2%	64.8%
	NaiveRAG	-Origin	NaiveRAG	-Origin	NaiveRAG	-Origin	NaiveRAG	-Origin
Comprehensiveness	24.8%	75.2%	39.2%	60.8%	16.4%	83.6%	44.4%	55.6%
Diversity	26.4%	<u>73.6%</u>	44.8%	<u>55.2%</u>	14.4%	<u>85.6%</u>	25.6%	<u>74.4%</u>
Empowerment	32.0%	68.0%	43.2%	<u>56.8%</u>	17.2%	82.8%	45.2%	<u>54.8%</u>
Overall	25.6%	<u>74.4%</u>	39.2%	60.8%	15.6%	84.4%	44.4%	<u>55.6%</u>

tailed insights. Consequently, this high-levelonly retrieval method may struggle with tasks that require precise, detailed answers.

• **Hybrid Mode**: The hybrid mode, or the full-version LightRAG, combines the strengths of both low- and high-level methods. It retrieves a broader set of relationships while simultaneously conducting an in-depth exploration of specific entities. This dual-level approach ensures both breadth in the retrieval process and depth in the analysis, providing a comprehensive view of the data. As a result, LightRAG achieves balanced performance across multiple dimensions.

Semantic Graph Excels in RAG. We eliminated the use of original text in our retrieval process. Surprisingly, the resulting variant, -Origin, does not exhibit significant performance declines. In some cases, this variant even shows improvements (e.g. Agriculture, Mix). We attribute this result to the effective extraction of key information during the graph-based indexing process, which provides sufficient context to answer queries. Additionally, the original text often contains irrelevant information that can introduce noise in the response.

4.4 Cost and Adaptability Analysis (RQ3)

We compare the cost of LightRAG with that of the top-performing baseline, GraphRAG, from two key perspectives. First, we examine the number of tokens and API calls during the indexing and

Table 3: RAG cost comparison on Legal data.

Phase	Retrieval Phase		Incremental Text Update	
Model	GraphRAG	Ours	GraphRAG	Ours
Tokens	$610 \times 1,000$	< 100	$1,399 \times 2 \times 5,000 \\ + T_{\text{extract}}$	$T_{ m extract}$
API Calls	$\frac{610\times1,000}{C_{\max}}$	1	$1,399 \times 2 + C_{\mathrm{extract}}$	$C_{ m extract}$

retrieval processes. Second, we analyze these metrics in relation to handling data changes in dynamic environments. The results on the legal dataset are presented in Table 3. In this context, $T_{\rm extract}$ represents the token overhead for entity and relationship extraction, $C_{\rm max}$ denotes the maximum number of tokens allowed per API call, and $C_{\rm extract}$ indicates the number of API calls required for extraction.

In the retrieval phase, GraphRAG generates 1,399 communities, with 610 level-2 communities actively utilized for retrieval in this experiment. Each community report averages 1,000 tokens, resulting in a total token consumption of 610,000 tokens (610 communities \times 1,000 tokens per community). Additionally, GraphRAG's requirement to traverse each community individually leads to hundreds of API calls, significantly increasing retrieval overhead. In contrast, LightRAG optimizes this process by using fewer than 100 tokens for keyword generation and retrieval, requiring only a single API call for the entire process. This efficiency is achieved through our retrieval mechanism, which seamlessly integrates graph structures and vectorized representations for information retrieval, thereby eliminating the need to process large volumes of information upfront.

In the incremental data update phase, addressing dynamic real-world scenarios, both models exhibit similar overhead for entity and relationship extraction. However, GraphRAG shows significant inefficiency managing newly added data. When introducing a new dataset matching the legal dataset size, GraphRAG must dismantle its existing community structure to incorporate new entities and relationships, then completely regenerate. This incurs substantial token costs of approximately 5,000 tokens per community report. With 1,399 communities, GraphRAG requires around 1,399 × 2 × 5,000 tokens to reconstruct both original and new community reports—an exorbitant expense highlighting its inefficiency. In contrast, LightRAG seamlessly integrates newly extracted entities and relationships into the existing graph without reconstruction, resulting in significantly lower overhead during incremental updates and demonstrating superior efficiency and cost-effectiveness.

4.5 Case Study

In addition to the overall evaluation results, we provided detailed case analyses in Appendix 9.2 to further illustrate the key findings of this study.

5 Related Work

5.1 Retrieval-Augmented Generation

Retrieval-Augmented Generation (RAG) systems enhance LLM inputs by retrieving relevant information from external sources, grounding responses in factual, domain-specific knowledge (Ram et al., 2023; Fan et al., 2024). Current RAG approaches (Gao et al., 2022, 2023; Chan et al., 2024; Yu et al., 2024) typically embed queries in a vector space to find the nearest context vectors. However, many of these methods rely on fragmented text chunks and only retrieve the top-k contexts, limiting their ability to capture comprehensive global information needed for high-quality responses.

Although recent studies (Edge et al., 2024) have explored using graph structures for knowledge representation, two key limitations persist. First, these approaches often lack the capability for dynamic updates and expansions of the knowledge graph, making it difficult to incorporate new information effectively. In contrast, our proposed model, LightRAG, addresses these challenges by enabling the RAG system to quickly adapt to new infor-

mation, ensuring the model's timeliness and accuracy. Additionally, existing methods often rely on brute-force searches for each generated community, which are inefficient for large-scale queries. Our LightRAG framework overcomes this limitation by facilitating rapid retrieval of relevant information from the graph through our proposed dual-level retrieval paradigm, significantly enhancing both retrieval efficiency and response speed.

5.2 Large Language Model for Graphs

Graphs have become a powerful framework for representing complex relationships, with applications across fields. As large language models (LLMs) evolve, researchers aim to enhance their ability to interpret graph-structured data. This work falls into three broad categories: i) GNNs as Prefix - using GNNs to generate structure-aware tokens for LLMs, e.g., GraphGPT (Tang et al., 2024) and LLaGA (Chen et al., 2024). ii) LLMs as Prefix leveraging LLMs to process graph data with text, producing embeddings/labels to refine GNN training, like GALM (Xie et al., 2023) and OFA (Liu et al., 2024). iii) LLMs-Graphs Integration achieving seamless interaction, using techniques,, fusion training, GNN alignment, and LLM-based agents for direct graph handling, e.g., Grenade (Li et al., 2023) and Congrat (Brannon et al., 2023).

6 Conclusion

This work introduces an advancement in LLMempowered Retrieval-Augmented Generation through the integration of a graph-based indexing approach that enhances both efficiency and comprehension in information retrieval. LightRAG utilizes a comprehensive knowledge graph to facilitate rapid and relevant document retrieval, enabling a deeper understanding of complex queries. Its dual-level retrieval paradigm allows for the extraction of both specific and abstract information, catering to diverse user needs. Furthermore, LightRAG's seamless incremental update capability ensures that the system remains current and responsive to new information, thereby maintaining its effectiveness over time. Overall, LightRAG excels in both efficiency and effectiveness, significantly improving the speed and quality of information retrieval and generation while reducing costs for LLM inference.

7 Limitations

Integrating multi-modal capabilities into LightRAG significantly enhances functionality by incorporating diverse data types such as text, images, and audio. This enables richer contextual understanding and more nuanced responses to user queries. Additionally, incorporating time-awareness ensures the system reflects dynamic events and evolving contexts, adapting responses based on temporal relevance. By combining multi-modal inputs with time-sensitive frameworks, LightRAG delivers timely, contextually accurate insights that address complex queries requiring understanding of current events and interconnections.

8 Acknowledgments

This research was partially supported by the National Natural Science Foundation of China under Grant No. U22B2019.

References

- William Brannon, Suyash Fulay, Hang Jiang, Wonjune Kang, Brandon Roy, Jad Kabbara, and Deb Roy. 2023. Congrat: Self-supervised contrastive pretraining for joint graph and text embeddings. *arXiv* preprint arXiv:2305.14321.
- Chi-Min Chan, Chunpu Xu, Ruibin Yuan, Hongyin Luo, Wei Xue, Yike Guo, and Jie Fu. 2024. Rq-rag: Learning to refine queries for retrieval augmented generation. *arXiv preprint arXiv:2404.00610*.
- Runjin Chen, Tong Zhao, AJAY KUMAR JAISWAL, Neil Shah, and Zhangyang Wang. 2024. Llaga: Large language and graph assistant. In *ICML*.
- Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, and Jonathan Larson. 2024. From local to global: A graph rag approach to query-focused summarization. arXiv preprint arXiv:2404.16130.
- Shahul Es, Jithin James, Luis Espinosa Anke, and Steven Schockaert. 2024. Ragas: Automated evaluation of retrieval augmented generation. In *EACL*.
- Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing Li. 2024. A survey on rag meeting llms: Towards retrieval-augmented large language models. In *KDD*.
- Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. 2022. Precise zero-shot dense retrieval without relevance labels. *arXiv preprint arXiv:2212.10496*.
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*.

- Yichuan Li, Kaize Ding, and Kyumin Lee. 2023. Grenade: Graph-centric language model for self-supervised representation learning on text-attributed graphs. In *EMNLP*.
- Hao Liu, Jiarui Feng, Lecheng Kong, Ningyue Liang, Dacheng Tao, Yixin Chen, and Muhan Zhang. 2024. One for all: Towards training one graph model for all classification tasks. In *ICLR*.
- Yuanjie Lyu, Zhiyu Li, Simin Niu, Feiyu Xiong, Bo Tang, Wenjin Wang, Hao Wu, Huanyong Liu, Tong Xu, and Enhong Chen. 2024. Crud-rag: A comprehensive chinese benchmark for retrieval-augmented generation of large language models. arXiv preprint arXiv:2401.17043.
- Hongjin Qian, Peitian Zhang, Zheng Liu, Kelong Mao, and Zhicheng Dou. 2024. Memorag: Moving towards next-gen rag via memory-inspired knowledge discovery. *arXiv preprint arXiv:2409.05591*.
- Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. In-context retrieval-augmented language models. *TACL*, 11:1316–1331.
- Ladislav Rampášek, Michael Galkin, Vijay Prakash Dwivedi, Anh Tuan Luu, Guy Wolf, and Dominique Beaini. 2022. Recipe for a general, powerful, scalable graph transformer. *NeurIPS*, 35:14501–14515.
- Alireza Salemi et al. 2024. Evaluating retrieval quality in retrieval-augmented generation. In *SIGIR*.
- Viju Sudhi, Sinchana Ramakanth Bhat, Max Rudat, et al. 2024. Rag-ex: A generic framework for explaining retrieval augmented generation. In *SIGIR*.
- Jiabin Tang, Yuhao Yang, Wei Wei, Lei Shi, Lixin Su, Suqi Cheng, Dawei Yin, and Chao Huang. 2024. Graphgpt: Graph instruction tuning for large language models. In *SIGIR*.
- Shangqing Tu, Yuanchun Wang, Jifan Yu, Yuyang Xie, Yaran Shi, Xiaozhi Wang, Jing Zhang, Lei Hou, and Juanzi Li. 2024. R-eval: A unified toolkit for evaluating domain knowledge of retrieval augmented large language models. In *KDD*.
- Han Xie, Da Zheng, Jun Ma, Houyu Zhang, Vassilis N Ioannidis, Xiang Song, Qing Ping, et al. 2023. Graphaware language model pre-training on a large graph corpus can help multiple graph applications. In *KDD*.
- Yue Yu, Wei Ping, Zihan Liu, Boxin Wang, Jiaxuan You, Chao Zhang, Mohammad Shoeybi, and Bryan Catanzaro. 2024. Rankrag: Unifying context ranking with retrieval-augmented generation in llms. *arXiv* preprint arXiv:2407.02485.
- Penghao Zhao, Hailin Zhang, Qinhan Yu, Zhengren Wang, Yunteng Geng, Fangcheng Fu, Ling Yang, Wentao Zhang, and Bin Cui. 2024. Retrieval-augmented generation for ai-generated content: A survey. arXiv preprint arXiv:2402.19473.

9 Appendix

In this section, we elaborate the methodologies and experimental settings of the LightRAG framework. It describes the specific steps for extracting entities and relationships from documents, detailing how large language models (LLMs) are utilized for this purpose. The section also specifies the prompt templates and configurations used in LLM operations, ensuring clarity in the experimental setup. Additionally, it outlines the evaluation criteria and dimensions used to assess the performance of LightRAG against baselines from various dimensions.

9.1 Experimental Settings

Table 4: Statistical information of the datasets.

Statistics	Agriculture	CS	Legal	Mix
# Docs	12	10	94	61
# Tokens	2,017,886	2,306,535	5,081,069	619,009

Evaluation Datasets. To conduct a comprehensive analysis of LightRAG, we selected four datasets from the UltraDomain benchmark (Qian et al., 2024). The UltraDomain data is sourced from 428 college textbooks and encompasses 18 distinct domains, including agriculture, social sciences, and humanities. Among these domains, we chose the Agriculture, CS, Legal, and Mix datasets. Each dataset contains between 600,000 and 5,000,000 tokens, with detailed information provided in Table 4. Below is a specific introduction to the four domains utilized in our experiments:

- **Agriculture**: This domain focuses on agricultural practices, covering a range of topics including beekeeping, hive management, crop production, and disease prevention.
- CS: This domain focuses on computer science and includes key areas of data science and software engineering. It particularly highlights machine learning and big data processing, featuring content on recommender systems, classification algorithms, and real-time analytics using Spark.
- Legal: It centers on corporate legal practices, addressing corporate restructuring, legal agreements, regulatory compliance, and governance, with a focus on the legal and financial sectors.
- **Mixed**: This domain presents a rich variety of literary, biographical, and philosophical texts, spanning a broad spectrum of disciplines, including cultural, historical, and philosophical studies.

Question Generation. To evaluate the effectiveness of RAG systems for high-level sensemaking tasks, we consolidate all text content from each dataset as context and adopt the generation method outlined in (Edge et al., 2024). Specifically, we instruct an LLM to generate five RAG users, along with five tasks for each user. Each generated user is accompanied by a textual description detailing their expertise and traits that motivate their question-raising activities. Each user task is also described, emphasizing one of the user's potential intentions when interacting with RAG systems. For each user-task pair, the LLM generates five questions that require an understanding of the entire corpus. In total, we generate 125 questions for each dataset.

Baselines. Our LightRAG is compared against the following state-of-the-art methods:

- Naive RAG (Gao et al., 2023): This model serves as a standard baseline in existing RAG systems. It segments raw texts into chunks and stores them in a vector database using text embeddings. For queries, Naive RAG generates vectorized representations to directly retrieve text chunks based on the highest similarity in their representations, ensuring efficient and straightforward matching.
- RQ-RAG (Chan et al., 2024): This approach leverages the LLM to decompose the input query into multiple sub-queries. These sub-queries are designed to enhance search accuracy by utilizing explicit techniques such as rewriting, decomposition, and disambiguation.
- HyDE (Gao et al., 2022): This method utilizes
 the LLM to generate a hypothetical document
 based on the input query. This generated document is then employed to retrieve relevant text
 chunks, which are subsequently used to formulate the final answer.
- GraphRAG (Edge et al., 2024): This is a graphenhanced RAG system that utilizes an LLM to extract entities and relationships from the text, representing them as nodes and edges. It generates corresponding descriptions for these elements, aggregates nodes into communities, and produces a community report to capture global information. When handling high-level queries, GraphRAG retrieves more comprehensive information by traversing these communities.

Implementation and Evaluation Details. In our experiments, we utilize the *nano vector database* for vector data management and access.

Table 5: Comparison of Document Insertion Times

No.	Token Count	LightRAG (s)	GraphRAG (s)
1	59,870	486	642
2	41,224	418	700
3	73,989	561	953
4	47,502	513	741
5	48,353	453	926

Table 6: Average Query Times

Metric	LightRAG	GraphRAG
Average Query Time (s)	11.2	23.6

For all LLM-based operations in LightRAG, we default to using GPT-4o-mini. To ensure consistency, the chunk size is set to 1200 across all datasets. Additionally, the gleaning parameter is fixed at 1 for both GraphRAG and LightRAG.

Defining ground truth for many RAG queries, particularly those involving complex high-level semantics, poses significant challenges. To address this, we build on existing work (Edge et al., 2024) and adopt an LLM-based multi-dimensional comparison method. We employ a robust LLM, specifically GPT-40-mini, to rank each baseline against our LightRAG. The evaluation prompt we used is detailed in Appendix 9.4.4. In total, we utilize four evaluation dimensions, including:

i) Comprehensiveness: How thoroughly does the answer address all aspects and details of the question? ii) Diversity: How varied and rich is the answer in offering different perspectives and insights related to the question? iii) Empowerment: How effectively does the answer enable the reader to understand the topic and make informed judgments? iv) Overall: This dimension assesses the cumulative performance across the three preceding criteria to identify the best overall answer.

The LLM directly compares two answers for each dimension and selects the superior response for each criterion. After identifying the winning answer for the three dimensions, the LLM combines the results to determine the overall better answer. To ensure a fair evaluation and mitigate the potential bias that could arise from the order in which the answers are presented in the prompt, we alternate the placement of each answer. We calculate win rates accordingly, leading to the final results.

9.2 Time and Space Comparison

To further investigate the scalability and efficiency of our proposed approach, LightRAG, we con-

Table 7: Final Storage Space Usage

Method	Final Storage Space (MB)
LightRAG	39.5
GraphRAG	286.7

ducted a series of incremental insertion experiments. We introduced five additional documents, with token counts ranging from 41,224 to 73,989, into the knowledge base. The results, presented in Table 5, demonstrate that LightRAG consistently outperforms the baseline GraphRAG method in both time and space efficiency. For the document indexing task, LightRAG exhibits near-linear scalability, with insertion times ranging from 418 to 561 seconds. In contrast, GraphRAG's insertion times are significantly higher, ranging from 642 to 953 seconds, indicating a heavier computational overhead due to its community detection mechanisms.

During the retrieval phase, LightRAG achieves an average query time of 11.2 seconds, less than half of GraphRAG's 23.6 seconds, as summarized in Table 6. This performance improvement is primarily attributed to LightRAG's lightweight, keyword-based retrieval approach. Furthermore, the final storage usage for LightRAG is only 39.5MB, as shown in Table 7, a stark contrast to GraphRAG's 286.7MB. This highlights LightRAG's superior efficiency in managing large-scale data. These results reinforce LightRAG's advantages in terms of scalability and resource efficiency, making it a more suitable choice for dynamic scenarios involving the ingestion and retrieval of large-scale knowledge bases.

9.3 Case Example of Retrieval-Augmented Generation in LightRAG

In Figure 2, we illustrate the retrieve-and-generate process. When presented with the query, "What metrics are most informative for evaluating movie recommendation systems?", the LLM first extracts both low-level and high-level keywords. These keywords guide the dual-level retrieval process on the generated knowledge graph, targeting relevant entities and relationships. The retrieved information is organized into three components: entities, relationships, and corresponding text chunks. This structured data is then fed into the LLM, enabling it to generate a comprehensive answer to the query.

9.4 Overview of Prompts in LightRAG

9.4.1 Prompts for Graph Generation

The graph construction prompt outlined in Figure 3 is designed to extract and structure entity-

```
Query: What metrics are most informative for evaluating movie recommendation systems?
                                                                                                                                                                         Querv
High level keywords: ["Metrics", "Movie recommendation systems", "Evaluation methods"],
Low level keywords: ["Accuracy", "Precision", "Recall", "F1 score", "User satisfaction", "Diversity", "Coverage"]
                                                                                                                                                                   Keywords
PERFORMANCE METRICS. Performance Metrics are measures used to evaluate how well the Decision Tree model performs, including RMSLE among others.

PRECISION-RECALL (PR) CURVE. The precision-recall curve is an event used to visualize the trade-off between precision and recall for different thresholds in binary
MOVIE CLUSTERING, Movie clustering refers to the grouping of movies based on similarities using the K-means algorithm and evaluating their performance metrics.
ACCURACY. Accuracy is a basic metric used to assess the performance of predictive models, expressed as a percentage of correct predictions
ITEMMATRIX. MOVIES. ItemMatrix holds the factors for movies, used to derive recommendations for each user based on their preferences.
AVGPRECISIONK, SPARK, avgPrecisionK function is utilized within the Spark framework to analyze the precision of recommendations for users.

MEAN SQUARE ERROR, VARIANCE, Mean Square Error is calculated based on the total Variance in the dataset, showing the relationship between explained and
ACTUALMOVIES, PREDICTEDMOVIES, ActualMovies are compared against predictedMovies to assess the accuracy of the movie recommendation model
       se metrics are popular for model evaluation and selection purposes ...
                                                                                  . These movie IDs are sorted in order of the estimated rating.\nAs we did for MSE and
RMSE, we can compute ranking-based evaluation metrics using MLlib's `RankingMetrics` class.
                                                                                                        ..... Then, for each user, we have the list of actual and predicted movie
IDs that we can pass to our APK function. In a manner similar to how we computed MSE, we will sum each of these APK scores using a 'reduce' action and divide by the
number of users (that is, the count of the `allRecs` RDD):\nFinally, we compute the accuracy for the decision tree...... First, we will compute the MSE and RMSE metri
Key Metrics for Evaluating Movie Recommendation Systems
Evaluating movie recommendation systems involves multiple metrics that assess the effectiveness and accuracy of the recommendations provided to users. Below are
some of the most informative metrics commonly used in this context, derived from the data and concepts surrounding recommendation systems
  Mean Average Precision at K (MAPK) [...]
2. Precision and Recall [...]
3. Root Mean Squared Error (RMSE) and Mean Squared Error (MSE) [...]
4. Area Under the Curve (AUC) [...]
5. F-Measure [...]
6. User Engagement Metrics [...]
Conclusion: A comprehensive evaluation of movie recommendation systems employs a range of metrics that capture both the accuracy of the recommendations and
user satisfaction. By effectively leveraging these metrics, system designers can optimize recommendation engines to deliver more personalized and relevant experiences
                                                                                                                                                           LLM Response
```

Figure 2: A retrieval and generation example.

relationship information from a text document based on specified entity types. The process begins by identifying entities and categorizing them into types such as organization, person, location, and event. It then provides detailed descriptions of their attributes and activities. Next, the prompt identifies relationships between these entities, offering explanations, assigning strength scores, and summarizing the relations using high-level keywords.

9.4.2 Prompts for Query Generation

In Figure 4, the query generation prompt outlines a framework for identifying potential user roles (e.g., data scientist, finance analyst, and product manager) and their objectives for generating queries based on a specified dataset description. The prompt explains how to define five distinct users who would benefit from interacting with the dataset. For each user, it specifies five key tasks they would perform while working with the dataset. Additionally, for each (user, task) combination, five high-level questions are posed to ensure a thorough understanding of the dataset.

9.4.3 Prompts for Keyword Extraction

In Figure 5, the prompt describes a method for extracting keywords from a user's query, distinguishing between high-level and low-level keywords. High-level keywords represent broad concepts or themes, while low-level keywords focus on specific entities and details. The extracted keywords are returned in JSON format, organized into two fields: "high_level_keywords" for overarching ideas and "low_level_keywords" for specific details.

9.4.4 Prompts for RAG Evaluation

The evaluation prompt is illustrated in Figure 6. It introduces a comprehensive evaluation framework for comparing two answers to the same question based on three key criteria: Comprehensiveness, Diversity, and Empowerment. Its purpose is to guide the LLM through the process of selecting the better answer for each criterion, followed by an overall assessment. For each of the three criteria, the LLM must identify which answer performs better and provide a rationale for its choice. Ultimately, an overall winner is determined based on

```
Given a text document that is potentially relevant to this activity and a list of entity types, identify all entities of those types from the text and all relationships among the
1. Identify all entities. For each identified entity, extract the following information:
 entity_name: Name of the entity, capitalized
- entity_type: One of the following types: [organization, person, geo. event]
            scription: Comprehensive description of the entity's attributes and activities
Format each entity as ("entity"<|><entity name><|><entity type><|><entity description>
2. From the entities identified in step 1, identify all pairs of (source_entity, target_entity) that are *clearly related* to each other.
For each pair of related entities, extract the following information:
- source_entity: name of the source entity, as identified in step 1
- target entity: name of the target entity, as identified in step 1
 - relationship_description: explanation as to why you think the source entity and the target entity are related to each other
- relationship_strength: a numeric score indicating strength of the relationship between the source entity and target entity
- relationship_keywords: one or more high-level key words that summarize the overarching nature of the relationship, focusing on concepts or themes rather than
Format each relationship as ("relationship"</><source_entity></scarget_entity></scretationship_description></scretationship_keywords></scretationship_strength>)
3. Identify high-level key words that summarize the main concepts, themes, or topics of the entire text. These should capture the overarching ideas present
Format the content-level key words as ("content_keywords"<|><high_level_keywords>)
4. Return output in English as a single list of all the entities and relationships identified in steps 1 and 2. Use **##** as the list delimiter.
5. When finished, output < |COMPLETE|>
-Real Data-
Entity_types: {entity_types}
Text: {input_text}
Output
                                                                                                                                      Graph Construct Prompt
```

Figure 3: Prompts for Graph Generation

```
Given the following description of a dataset: {total_description}

Please identify 5 potential users who would engage with this dataset. For each user, list 5 tasks they would perform with this dataset. Then, for each (user, task) combination, generate 5 questions that require a high-level understanding of the entire dataset.

Output the results in the following structure:

- User 1: [user description]

- Task 1: [task description] [ Question 1: {Question 1}, Question 2: {Question 2}, Question 3: {Question 3}, Question 4: {Question 4}, Question 5: {Question 5} ]

- Task 2: [task description] [ Question 1: {Question 1}, Question 2: {Question 2}, Question 3: {Question 3}, Question 4: {Question 4}, Question 5: {Question 5} ]

- Task 5: [task description] [ Question 1: {Question 1}, Question 2: {Question 2}, Question 3: {Question 3}, Question 4: {Question 4}, Question 5: {Question 5} ]

- User 2: [user description]

- User 5: [user description]
```

Figure 4: Prompts for Query Generation

performance across all three dimensions, accompanied by a detailed summary that justifies the decision. The evaluation is structured in JSON format, ensuring clarity and consistency, and facilitating a systematic comparison between the two answers.

9.5 Case Study: Comparison between LightRAG and the Baseline NaiveRAG

To further illustrate LightRAG's superiority over baseline models in terms of comprehensiveness, empowerment, and diversity, we present a case study comparing LightRAG and NaiveRAG in Table 8. This study addresses a question regarding indigenous perspectives in the context of corporate mergers. Notably, LightRAG offers a more in-depth exploration of key themes related to indigenous perspectives, such as cultural significance, collaboration, and legal frameworks, supported by specific and illustrative examples. In contrast,

while NaiveRAG provides informative responses, it lacks the depth needed to thoroughly examine the various dimensions of indigenous ownership and collaboration. The dual-level retrieval process employed by LightRAG enables a more comprehensive investigation of specific entities and their interrelationships, facilitating extensive searches that effectively capture overarching themes and complexities within the topic.

9.6 Case Study: Comparison between LightRAG and GraphRAG

To provide a clear comparison between baseline methods and our LightRAG, we present specific case examples in Table 9. The table includes responses to a machine learning question from both the competitive baseline, GraphRAG, and our LightRAG framework. In this instance, LightRAG outperforms GraphRAG in all evaluation dimensions

```
You are a helpful assistant tasked with identifying both high-level and low-level keywords in the user's query.
Given the query, list both high-level and low-level keywords. High-level keywords focus on overarching concepts or themes, while low-level keywords focus on
specific entities, details, or concrete terms.
                                                                                                                      Keywords Generate Instruction Prompt

    Output the keywords in JSON format.
    The JSON should have two keys:

- "high_level_keywords" for overarching concepts or themes.
- "low_level_keywords" for specific entities or details.
Example 1: Query: "How does international trade influence global economic stability?
Output: {{ "high_level_keywords": ["International trade", "Global economic stability", "Economic impact"], "low_level_keywords": ["Trade agreements", "Tariffs", "Currency exchange", "Imports", "Exports"] }}
Query: "What are the environmental consequences of deforestation on biodiversity?"
Output: {{ "high_level_keywords": ["Environmental consequences", "Deforestation", "Biodiversity loss"], "low_level_keywords": ["Species extinction", "Habitat destruction", "Carbon emissions", "Rainforest", "Ecosystem"] }}
        "What is the role of education in reducing poverty?"
Output: {{ "high_level_keywords": ["Education", "Poverty reduction", "Socioeconomic development"], "low_level_keywords": ["School access", "Literacy rates", "Job
training", "Income inequality"] }}
Query: {query}
Output:
                                                                                                                                 Keywords Generate Input Prompt
```

Figure 5: Prompts for Keyword Extraction

```
---Role---
You are an expert tasked with evaluating two answers to the same question based on four criteria: Comprehensiveness, Diversity, and Empowerment.

---Goal---
You will evaluate two answers to the same question based on four criteria: Comprehensiveness, Diversity, and Empowerment.

- Comprehensiveness: How much detail does the answer provide to cover all aspects and details of the question?
- Diversity: How varied and rich is the answer in providing different perspectives and insights on the question?
- Empowerment: How well does the answer help the reader understand and make informed judgments about the topic?

For each criterion, choose the better answer (either Answer 1 or Answer 2) and explain why. Then, select an overall winner based on these three categories.

Evaluation Instruction Prompt

Here is the question: {query}

Here are the two answers: Answer 1: {answer1}; Answer 2: {answer2}

Evaluate both answers using the three criteria listed above and provide detailed explanations for each criterion.

Output your evaluation in the following JSON format:

("Comprehensiveness": {("Winner": "[Answer 1 or Answer 2]", "Explanation": "[Provide explanation here]" }},

"Diversity": {{ "Winner": "[Answer 1 or Answer 2]", "Explanation": "[Provide explanation here]" }},

"Overall Winner": "[Answer 1 or Answer 2]", "Explanation": "[Summarize why this answer is the overall winner based on the three criteria]" }}
```

Figure 6: Prompts for RAG Evaluation

assessed by the LLM judge, including comprehensiveness, diversity, empowerment, and overall quality. Our key observations are as follows:

- i) Comprehensiveness. Notably, LightRAG covers a broader range of machine learning metrics, showcasing its comprehensiveness and ability to effectively discover relevant information. This highlights the strength of our graph-based indexing paradigm, which excels in precise entity and relation extraction as well as LLM profiling.
- ii) Both Diversity and Empowerment. Furthermore, LightRAG not only offers a more diverse array of information but also delivers more empowering content. This success is due to LightRAG's

innovative hierarchical retrieval paradigm. This approach combines in-depth explorations of related entities through low-level retrieval to enhance the empowerment and depth of the responses, along with broader explorations via high-level retrieval to improve the diversity of the answers. Together, they capture a comprehensive global perspective of the knowledge domain, enabling LightRAG to provide users with a richer, more empowering, and more diverse set of information.

Evaluation Input Prompt

Table 8: Case Study: Comparison Between LightRAG and the Baseline NaiveRAG.

Query	How do indigenous perspectives on ownership and collaboration influence corporate mergers in countries like Canada and Australia?
NaiveRAG	Indigenous perspectives on ownership and collaboration play a significant role in shaping corporate mergers in countries like Canada and Australia, particularly where resources and land right intersect with Indigenous sovereignty and cultural values. These perspectives emphasize community engagement, respect for traditional land use, and a collaborative approach to resource management influencing corporate strategies and outcomes.
	Both nations illustrate that prioritizing Indigenous perspectives leads to more sustainable and equitable outcomes. Acknowledging Indigenous rights and fostering collaborative relationships can reduce conflicts and enhance corporate reputations. Ultimately, integrating Indigenous perspective is a pathway to better business practices and community relations.
LightRAG	Indigenous Perspectives on Ownership and Collaboration in Corporate Mergers Indigenous perspectives significantly influence corporate mergers in countries like Canada and
	Australia. These views emphasize communal rights to land and resources, spiritual connections to the environment, and the importance of collaboration over competition.
	Cultural Significance of Land Ownership [] The Role of Collaboration []
	Legal and Regulatory Frameworks []
	Conclusion
LLM Decision	Indigenous perspectives shape corporate mergers by emphasizing communal ownership and collab orative partnerships. Corporations must adopt respectful, inclusive practices that honor Indigenou cultural values and rights, paving the way for sustainable, mutually beneficial outcomes. Comprehensiveness: "Winner": "Answer 2 (LightRAG)",
	"Explanation": "Answer 2 provides a thorough exploration of the indigenous perspective by discussing cultural significance, collaboration, and legal frameworks, including specific examples and detailed insights. In contrast, Answer 1, while informative, lacks the same depth in analyzing the various dimensions of the indigenous perspectives on ownership and collaboration.
	Diversity : "Winner": "Answer 2 (LightRAG)", "Explanation": "Answer 2 presents a wider array of perspectives, focusing on the communal aspector of land ownership, spiritual connections, and practical examples of collaboration. It covers more ground by contrasting Indigenous views with Western notions, whereas Answer 1 primarily focuse on corporate strategies and outcomes with limited perspective."
	Empowerment : "Winner": "Answer 2 (LightRAG)", "Explanation": "Answer 2 equips the reader with nuanced understanding and actionable insights by highlighting the importance of collaboration and the legal frameworks involved. I empowers corporations by illustrating how adopting an inclusive approach can facilitate better ou comes. Answer 1, while practical, does not emphasize the moral or ethical implications as strongly.
	Overall Winner: "Winner": "Answer 2 (LightRAG)", "Explanation": "Answer 2 excels overall due to its comprehensive exploration, diversity of perspectives, and empowerment of the reader with actionable insights about indigenous perspectives and collaboration in corporate mergers. Although Answer 1 is more direct, the depth and breadth of Answer 2 make it the stronger response."

Table 9: Case Study: Comparison Between LightRAG and the Baseline Method GraphRAG.

Query	Which methods can normalize the feature values for effective machine learning?
GraphRAG	Key Metrics for Evaluating Movie Recommendation Systems: Evaluating the effectiveness
	of movie recommendation systems is crucial for ensuring user satisfaction and improving the
	overall performance of these systems. Several metrics are widely recognized for their ability to
	provide insights into how well a recommendation system performs. Below are the most informative
	metrics: 1. Precision [] 2. Recall [] 3. F1 Score [] 4. Mean Average Precision (MAP)
	[] 5. Root Mean Squared Error (RMSE) [] 6. User Satisfaction Metrics [] Conclusion:
	In summary, metrics such as precision, recall, F1 score, MAP, RMSE, and user satisfaction are
	critical for evaluating the performance of movie recommendation systems. Each metric provides unique insights that can help developers and researchers refine their algorithms and enhance user
	experiences. By focusing on these metrics, stakeholders can ensure that their recommendation
	systems effectively meet user preferences and expectations.
LightRAG	Key Metrics for Evaluating Movie Recommendation Systems: Evaluating movie recommen-
g	dation systems involves multiple metrics that assess the effectiveness and accuracy of the recom-
	mendations provided to users. Below are some of the most informative metrics commonly used in
	this context, derived from the data and concepts surrounding recommendation systems. 1. Mean
	Average Precision at K (MAPK) [] 2. Precision and Recall [] 3. Root Mean Squared
	Error (RMSE) and Mean Squared Error (MSE) [] 4. Area Under the Curve (AUC) [] 5.
	F-Measure [] 6. User Engagement Metrics [] Conclusion : A comprehensive evaluation of
	movie recommendation systems employs a range of metrics that capture both the accuracy of the
	recommendations and user satisfaction. By effectively leveraging these metrics, system designers
IIMP	can optimize recommendation engines to deliver more personalized and relevant experiences.
LLM Decision	Comprehensiveness: "Winner": "Answer 2 (LightRAG)",
	"Explanation": "Answer 2 provides a broader array of metrics including MAPK, AUC, and user engagement metrics alongside precision, recall, RMSE, and F-Measure. This gives a fuller
	view of how recommendation systems can be evaluated. Answer 1, while detailed, is more
	focused on traditional metrics and lacks the inclusion of engagement metrics which are crucial for
	understanding user interaction."
	Diversity : "Winner": "Answer 2 (LightRAG)",
	"Explanation": "Answer 2 not only covers a wide variety of metrics but also includes nuanced
	explanations of how some metrics interrelate and differ from one another, like the inclusion of
	both RMSE and MSE, as well as the details behind AUC. In contrast, Answer 1 sticks primarily to
	standard metrics without much exploration of potential nuances."
	Empowerment: "Winner": "Answer 2 (LightRAG)",
	"Explanation": "Answer 2 empowers the reader more effectively by detailing how each metric
	functions and its importance in evaluating recommendation systems. By providing context such as
	the trade-offs between precision and recall and emphasizing user engagement metrics, it enables
	readers to make more informed judgments and understand the implications of different metrics.
	Answer 1 is more straightforward but lacks the depth of insight regarding why these metrics matter."
	Overall Winner: "Winner": "Answer 2 (LightRAG)",
	"Explanation": "While Answer 1 is more direct and systematic, Answer 2 excels in comprehen-
	siveness, diversity, and empowerment. It provides a richer exploration of the topic, including
	insights into user engagement and nuanced differences between metrics. This depth and breadth
	make it more informative for readers seeking to thoroughly understand the evaluation of movie
	recommendation systems."