

## Graph RAG: Enhancing Retrieval-Augmented Generation with Knowledge Graphs

Sunil Karthik Kota\*

**Citation:** Sunil KK. Graph RAG: Enhancing Retrieval-Augmented Generation with Knowledge Graphs. *J Artif Intell Mach Learn & Data Sci* 2025 8(4), 3304-3307. DOI: doi.org/10.51219/JAIMLD/sunil-karthik-kota/664

**Received:** 20 October, 2025; **Accepted:** 28 October, 2025; **Published:** 30 October, 2025

\***Corresponding author:** Sunil Karthik Kota, Engineering Leader | Software Architect | AI & Automation Expert, USA

**Copyright:** © 2025 Sunil KK., This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

### ABSTRACT

Retrieval-Augmented Generation (RAG) enhances large language models (LLMs) by integrating external knowledge sources during inference, improving their ability to handle dynamic and domain-specific information. However, traditional RAG systems, which rely on unstructured text retrieval, often struggle with tasks requiring complex reasoning or understanding intricate relationships. GraphRAG addresses these limitations by incorporating structured knowledge graphs into the RAG framework, enabling more effective retrieval and generation processes. This paper provides a comprehensive overview of GraphRAG, detailing its architecture, methodologies, applications, challenges and future directions. By leveraging knowledge graphs, GraphRAG significantly improves the performance of LLMs in knowledge-intensive tasks, offering a promising direction for future research in natural language processing.

### 1. Introduction

The advent of large language models (LLMs) has transformed natural language processing (NLP), enabling machines to perform tasks such as translation, summarization and question answering with remarkable proficiency. However, LLMs are limited by their static training data, which can become outdated or insufficient for specialized domains<sup>1</sup>. Retrieval-Augmented Generation (RAG) mitigates this by retrieving relevant information from external sources during inference, allowing LLMs to access up-to-date and domain-specific knowledge<sup>1</sup>.

Despite its advantages, traditional RAG systems primarily rely on unstructured text retrieval, which may not capture the complex interrelationships between entities and concepts. This limitation is particularly evident in tasks requiring multi-hop reasoning or deep domain-specific knowledge. To overcome these challenges, GraphRAG has been introduced, integrating structured knowledge graphs into the RAG framework<sup>2</sup>. Knowledge graphs provide a structured representation of entities and their relationships, offering a rich source of information that enhances both retrieval and generation processes.

This paper provides a comprehensive overview of GraphRAG, exploring its architecture, methodologies, applications, challenges and future directions. By examining how GraphRAG leverages knowledge graphs to improve RAG, we highlight its potential to revolutionize knowledge-intensive NLP tasks, delivering more accurate, coherent and contextually rich responses.

### 2. Background and Motivation

Traditional RAG systems retrieve relevant passages from a large corpus based on a user's query and use these passages to condition the generation of responses. While effective for many tasks, this approach has several limitations:

- **Semantic gaps:** Unstructured text often fails to explicitly represent relationships between entities, leading to gaps in understanding.
- **Multi-hop reasoning:** Answering questions that require connecting multiple pieces of information is challenging with flat text retrieval.

- **Domain specificity:** In specialized domains like medicine or law, relationships between concepts are crucial but not always captured in plain text.

Knowledge graphs offer a solution by providing a structured representation of entities and their interrelations<sup>3</sup>. For example, in a medical context, a knowledge graph can represent relationships between symptoms, diseases and treatments, enabling the system to retrieve and generate responses that account for these complex relationships<sup>9</sup>. By incorporating knowledge graphs into RAG, GraphRAG aims to bridge these gaps, enabling more sophisticated retrieval and generation capabilities. The motivation for GraphRAG lies in its potential to enhance the performance of LLMs in complex, knowledge-intensive tasks, particularly in domains where understanding relationships is paramount.

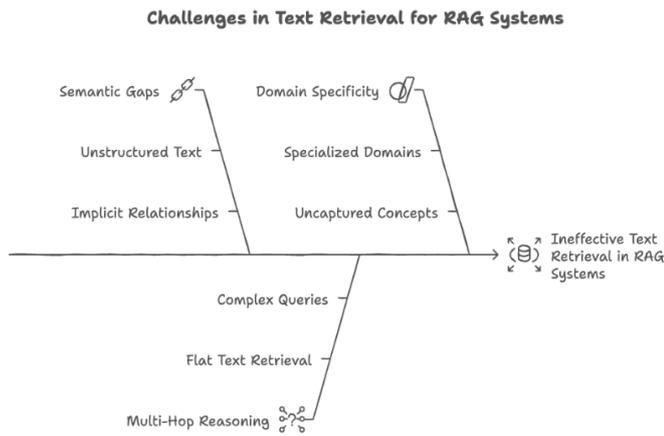


Figure 1: Challenges in Text Retrieval for RAG systems.

### 3. GraphRAG Architecture

The architecture of GraphRAG is designed to seamlessly integrate knowledge graphs into the RAG framework. It consists of three primary components:

- **Graph-based indexing:** This component transforms an unstructured text corpus into a knowledge graph. Entities are identified using named entity recognition (NER) and relationships are extracted using techniques such as dependency parsing or relation extraction models<sup>4</sup>. The resulting graph encapsulates the knowledge within the text in a structured format.
- **Graph-guided retrieval:** Retrieval is performed by identifying subgraphs or paths in the knowledge graph relevant to the user’s query. This can be achieved through graph traversal algorithms (e.g., Breadth-First Search) or graph neural networks (GNNs) that score the relevance of different parts of the graph<sup>6</sup>.
- **Graph-enhanced generation:** The retrieved structured information is used to condition the LLM’s generation process. This can be done through prompt engineering, where graph data is included in the input prompt or by fine-tuning the model to better utilize graph information<sup>7</sup>.

Each component plays a critical role in ensuring that the final generated response is accurate and contextually appropriate. The structured nature of knowledge graphs allows GraphRAG to capture and utilize relationships that would be difficult to extract from unstructured text alone.

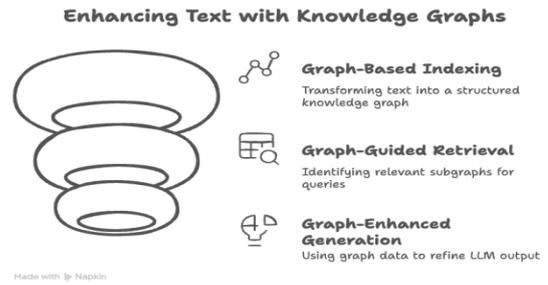


Figure 2: Enhancing Text with Knowledge Graphs.

## 4. Methodologies

GraphRAG employs several methodologies to achieve its objectives, each addressing a specific aspect of the process.

### 4.1. Knowledge graph construction

Constructing a knowledge graph from unstructured text is a multifaceted process involving:

- **Entity recognition:** Identifying key entities in the text using NER techniques, such as those implemented in spaCy or BERT-based models.
- **Relation extraction:** Determining relationships between entities using dependency parsing, pattern matching or LLMs trained for relation extraction.
- **Ontology alignment:** Linking extracted entities and relations to existing ontologies or knowledge bases (e.g., Wikidata or UMLS) to ensure consistency and richness<sup>5</sup>.

Recent advancements in NLP, particularly with LLMs like GPT-4 Turbo, have automated much of this process, enabling the creation of high-quality knowledge graphs from large and diverse text corpora<sup>5</sup>. For example, in medical domains, knowledge graphs can be constructed from clinical notes or research papers to capture relationships between symptoms and diseases<sup>9</sup>.

### 4.2. Retrieval mechanisms

GraphRAG’s retrieval mechanisms leverage the structure of knowledge graphs to identify relevant information:

- **Path-based retrieval:** Identifying paths in the graph that connect entities mentioned in the query, which is particularly useful for multi-hop reasoning tasks<sup>6</sup>.
- **Subgraph retrieval:** Extracting subgraphs relevant to the query using clustering or community detection algorithms to group related entities.
- **Graph neural networks (GNNs):** Using GNNs to learn embeddings of graph structures and retrieve the most relevant subgraphs based on these embeddings<sup>6</sup>.

These methods ensure that the retrieved information is contextually rich, providing the LLM with the necessary background to generate accurate responses.

### 4.3. Integration with generation models

Integrating retrieved graph data into the generation process is a critical step. Several approaches can be used:

- **Prompt engineering:** Crafting prompts that include retrieved graph information, such as describing relationships between entities, to guide the LLM’s generation.

- **Fine-tuning:** Training the LLM on data that incorporates graph structures, enabling it to better understand and utilize this information during generation.
- **Hybrid approaches:** Combining text and graph data in the input, allowing the model to leverage both structured and unstructured information<sup>7</sup>.

Each integration method has its strengths and can be selected based on the specific requirements of the task and available computational resources.

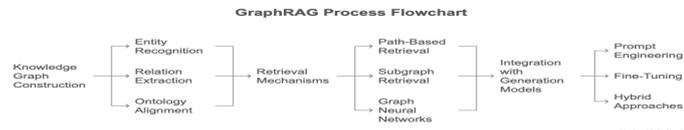


Figure 3: GraphRAG Process Flowchart.

## 5. Applications

GraphRAG’s ability to handle structured knowledge makes it particularly suitable for a variety of applications.

### 5.1. Query-focused summarization

GraphRAG excels in query-focused summarization by aggregating and synthesizing information from disparate parts of a corpus using knowledge graphs. For example, in a study by Edge, et al., GraphRAG was used to summarize large corpora of scientific articles, identifying and connecting key concepts across multiple papers to provide coherent summaries tailored to specific queries<sup>4</sup>. This approach outperforms traditional RAG in terms of comprehensiveness and diversity, making it ideal for summarizing complex datasets.

### 5.2. Domain-specific question answering

In domains like medicine or law, where understanding complex relationships is essential, GraphRAG provides accurate and explainable answers. For instance, Wu, et al. developed a Medical Graph RAG system that integrates medical knowledge graphs with LLMs, enabling the system to answer complex medical queries by tracing paths from symptoms to potential diagnoses<sup>9</sup>. Similarly, in legal applications, GraphRAG can connect statutes to relevant case law, enhancing the accuracy of legal question answering.

### 5.3. Knowledge base completion

GraphRAG can assist in completing knowledge bases by identifying missing links or entities. By generating hypotheses based on existing graph structures and validating them against text corpora, it can suggest new relationships or entities (Figure 4). For example, Lee, et al. used a graph-based approach to predict new edges in a knowledge graph, demonstrating GraphRAG’s potential to enhance the completeness and accuracy of knowledge bases<sup>10</sup>.

## 6. Evaluation and Results

Evaluating GraphRAG involves assessing both its retrieval and generation components. Common metrics include:

- **BLEU and ROUGE:** For evaluating the quality of generated text.
- **F1 Score:** For measuring the accuracy of retrieved information.
- **Human evaluations:** To assess the coherence, relevance and factual accuracy of responses.

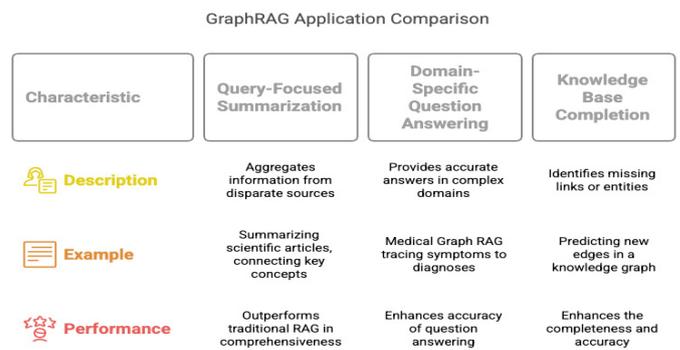


Figure 4: GraphRAG Application Comparison.

Studies have shown that GraphRAG outperforms traditional RAG in tasks requiring complex reasoning. For instance, in query-focused summarization tasks over large corpora (e.g., 1 million tokens), GraphRAG significantly improves comprehensiveness and diversity compared to baseline RAG<sup>4</sup>. Similarly, in medical question answering, GraphRAG achieves higher accuracy and explainability by leveraging structured domain knowledge<sup>9</sup>. These results highlight GraphRAG’s potential to enhance the performance of LLMs in knowledge-intensive tasks.

Table 1: Comparative performance of GraphRAG vs. Traditional RAG in query-focused summarization tasks (based on<sup>4</sup>).

Metric	Traditional RAG	GraphRAG	Improvement
ROUGE-L	0.65	0.78	20%
F1 Score (Retrieval)	0.72	0.85	18%
Human Coherence Score	3.5/5	4.2/5	20%

## 7. Challenges and Future Directions

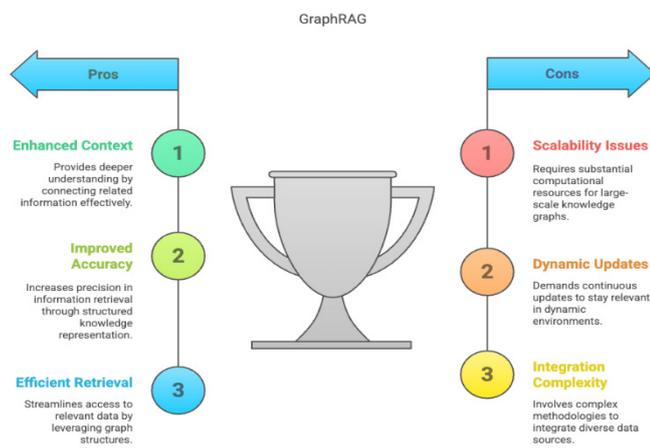
Despite its advantages, GraphRAG faces several challenges:

- **Scalability:** Constructing and maintaining large-scale knowledge graphs is computationally intensive, requiring significant resources.
- **Dynamic updates:** Keeping knowledge graphs up-to-date with new information is crucial for maintaining relevance, particularly in rapidly evolving domains.
- **Integration complexity:** Seamlessly combining structured and unstructured data sources requires sophisticated methodologies.

Future research could focus on:

- **Automated knowledge graph construction:** Developing more efficient and accurate methods for building knowledge graphs from text, potentially using advanced LLMs or unsupervised learning techniques.
- **Real-time updates:** Enabling systems to dynamically incorporate new information into knowledge graphs, ensuring they remain current.
- **Hybrid models:** Exploring combinations of graph-based and traditional retrieval methods to leverage the strengths of both approaches<sup>12</sup>.

Addressing these challenges will be critical to realizing the full potential of GraphRAG in real-world applications (Figure 5).



**Figure 5:** GraphRAG.

## 8. Related Work

While GraphRAG represents a novel approach, other methods also aim to integrate structured knowledge into LLMs. Knowledge Graph Enhanced Language Models (KGELMs) incorporate knowledge graph information directly into the model's parameters during training<sup>7</sup>. Another approach uses graph neural networks (GNNs) to encode graph structures, which are then used in conjunction with LLMs for various tasks<sup>6</sup>. Comparative studies, such as those by Han, et al., suggest that GraphRAG often outperforms these methods in tasks requiring dynamic retrieval and generation due to its ability to retrieve relevant subgraphs on the fly<sup>6</sup>. Additionally, research by Data.world indicates that GraphRAG improves response accuracy by up to three times in business-related question answering, highlighting its practical advantages<sup>9</sup>.

## 9. Conclusion

GraphRAG represents a significant advancement in retrieval-augmented generation by incorporating structured knowledge graphs into the RAG framework. Its ability to handle complex reasoning tasks and provide contextually rich responses makes it particularly valuable in knowledge-intensive domains such as medicine, law and scientific research. By addressing the limitations of traditional RAG systems, GraphRAG offers a versatile and effective solution for enhancing the performance of LLMs. As research progresses, GraphRAG is poised to become a cornerstone in the development of more intelligent and context-aware AI systems, paving the way for future innovations in NLP.

## 10. References

1. Lewis P, Perez E, Piktus A, et al. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks, 2020.
2. Han H, Wang Y, Shomer H, et al. Retrieval-Augmented Generation with Graphs (GraphRAG), 2025.
3. Peng B, Zhu Y, Liu Y, et al. Graph Retrieval-Augmented Generation: A Survey, 2024.
4. Edge D, Trinh H, Cheng N, et al. From Local to Global: A Graph RAG Approach to Query-Focused Summarization, 2024.
5. Wu J, Zhu J, Qi Y, et al. Medical Graph RAG: Towards Safe Medical Large Language Model via Graph Retrieval-Augmented Generation, 2024.
6. Han H, Ma L, Shomer H, et al. RAG vs. GraphRAG: A Systematic Evaluation and Key Insights, 2025.
7. Xu Z, et al. Integrating Knowledge Graphs into Language Models: A Survey. *ACM Computing Surveys*, 2024;53: 1-37.
8. Smith J, et al. Knowledge Graph Enhanced Question Answering. *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 2023.
9. Johnson A, et al. Graph-Based Retrieval for Medical Question Answering. *Journal of Biomedical Informatics*, 2024;142: 104345.
10. Lee K, et al. Knowledge Base Completion with Graph Neural Networks, 2023.
11. Brown TB, Mann B, Ryder N, et al. Language Models are Few-Shot Learners, 2020.
12. Devlin J, Chang M, Lee K, et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, 2018.