

GRAPH RETRIEVAL AUGMENTED GENERATION

Yogesh Haribhau Kulkarni



About Me

YHK

Yogesh Haribhau Kulkarni

Bio:

- ▶ 20+ years in CAD/Engineering software development
- ▶ Got Bachelors, Masters and Doctoral degrees in Mechanical Engineering (specialization: Geometric Modeling Algorithms).
- ▶ Currently doing Coaching in fields such as Data Science, Artificial Intelligence Machine-Deep Learning (ML/DL) and Natural Language Processing (NLP).
- ▶ Feel free to follow me at:
 - ▶ Github (github.com/yogeshhk)
 - ▶ LinkedIn (www.linkedin.com/in/yogeshkulkarni/)
 - ▶ Medium (yogeshharibhaukulkarni.medium.com)
 - ▶ Send email to [yogeshkulkarni at yahoo dot com](mailto:yogeshkulkarni@yahoo.com)



Office Hours:
Saturdays, 2 to 5pm
(IST); Free-Open to all;
email for appointment.

Introduction to Graph RAG

YHK

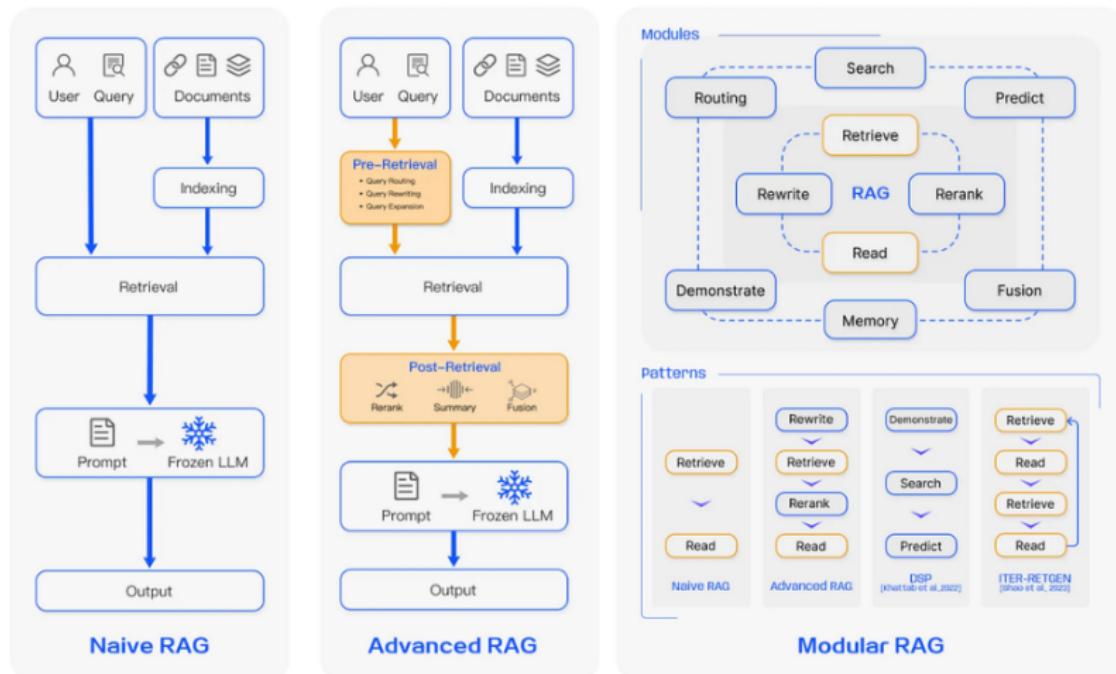
Challenges with LLMs

- ▶ Learns random sentences from random people
- ▶ Talks like a person but doesn't really understand what it's saying
- ▶ Occasionally speaks absolute non sense
- ▶ Sensitive to question phrasing
- ▶ Limited to public "knowledge"

(Ref: The GenAI Stack - Andreas Kollegger - Neo4j)



Thematic RAG Classification



(Ref: <https://graphrag.com/concepts/intro-to-graphrag/>)

The phases of an advanced RAG system

- ▶ Pre-retrieval-Query rewriting, query entity extraction, query expansion, etc.
- ▶ Retrieval of relevant context
- ▶ Post-retrieval: Reranking, pruning, etc.
- ▶ Answer generation

(Ref: <https://graphrag.com/concepts/intro-to-graphrag/>)

Why Graph RAG?

- ▶ Language models struggle with factual accuracy and real-world knowledge.
- ▶ Retrieval-Augmented Generation (RAG) improves accuracy using external text data.
- ▶ Traditional RAG has limitations in context understanding and scalability.
- ▶ GraphRAG leverages knowledge graphs for better retrieval and response generation.



Limitations of Traditional RAG

- ▶ **Flat Retrieval:** Documents are treated as isolated entities.
- ▶ **Contextual Shortcomings:** Lacks deep semantic understanding.
- ▶ **Scalability Issues:** Slower retrieval with increasing data volume.
- ▶ Struggles with scalability due to flat data representation.
- ▶ Lacks global context over the entire data corpus.
- ▶ Inefficient for complex reasoning across multiple documents.
- ▶ **Lack of Explainability:** Hard to trace the source of retrieved information.
- ▶ **Local Window:** Limited chunk-level context leads to incomplete responses.
- ▶ **Loss of Structural Relationships:** Ignores hierarchical relationships in data.

Challenges with Microsoft's GraphRAG

- ▶ Too expensive and impractical for large-scale industrial use.
- ▶ Most companies prefer standard vector databases.
- ▶ Lack of widespread production adoption due to cost and complexity.



Text-To-Cypher/SPARQL Alternative

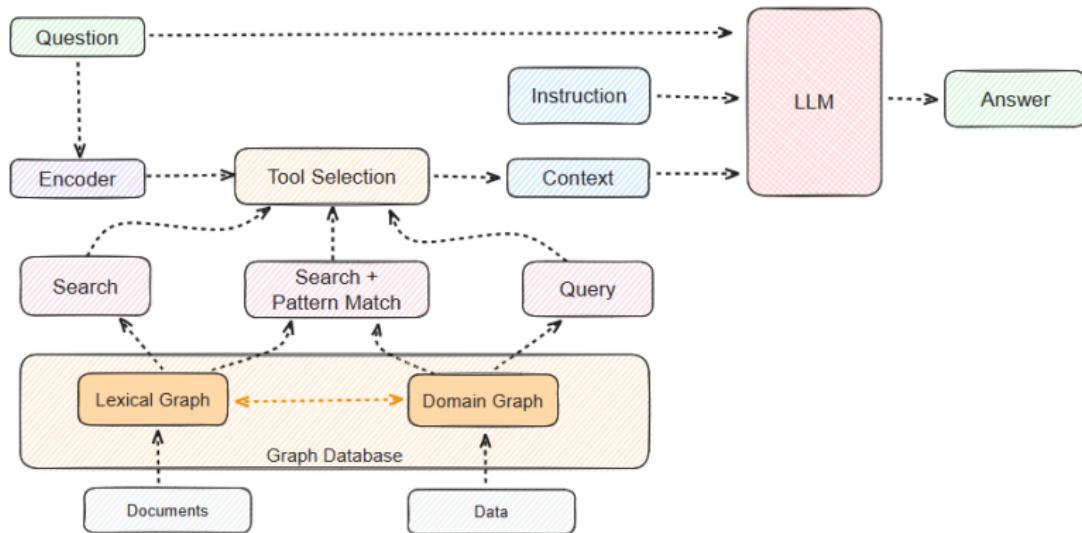
- ▶ Effective alternative to Microsoft's GraphRAG.
- ▶ Requires costly LLM calls for query generation.
- ▶ Adds uncertainty due to prompt effectiveness and model performance.
- ▶ Increases response time and implementation complexity.



What is GraphRAG?

- ▶ No universally accepted definition yet.
- ▶ Some associate it with Microsoft's graph-based search approach.
- ▶ Others define it as querying LPG or RDF graphs using LLM-generated queries (Cypher, SPARQL).
- ▶ Uses knowledge graphs instead of unstructured text.
- ▶ Captures entities, relationships, and hierarchical structures.
- ▶ Enables accurate, context-aware retrieval and response generation.
- ▶ Supports complex query handling with enhanced explainability.
- ▶ Combines structured Knowledge Graphs (KGs) with semantic vectors.
- ▶ Enables LLMs to reason over multi-hop connections.
- ▶ Provides a holistic perspective by connecting different data sources.

Basics of GraphRAG



(Ref: [https://graphrag.com/concepts/intro-to-graphrag/\(\)](https://graphrag.com/concepts/intro-to-graphrag/))

Why RAG?

- ▶ RAG is widely used for real-world enterprise applications.
- ▶ Augments external knowledge sources to query private corpora.
- ▶ Traditional vector-based RAG relies only on semantic similarity.

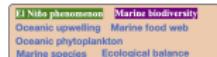
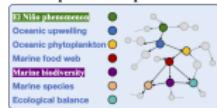
Key Features of GraphRAG

- ▶ **Structured Representation:** Uses knowledge graphs.
- ▶ **Contextual Retrieval:** Understands semantic relationships.
- ▶ **Efficient Processing:** Reduces computational cost.
- ▶ **Multi-Faceted Queries:** Synthesizes data from multiple sources.
- ▶ **Explainability:** More transparent than black-box models.
- ▶ **Continuous Learning:** Expands knowledge over time.



Query:

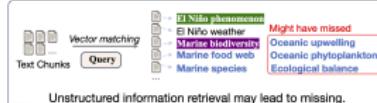
"How does the El Niño phenomenon potentially impact marine biodiversity?"

Related Entities:**Multi-hop Relationships:****Multi-hop Reasoning:**

- El Niño phenomenon impacts oceanic upwelling.
- Changes in upwelling affect the levels of oceanic phytoplankton.
- Oceanic phytoplankton levels influence marine biodiversity.

Answer:

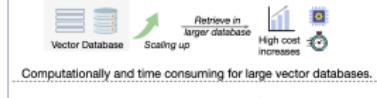
"The El Niño phenomenon reduces nutrient-rich upwelling, leading to a decline in oceanic phytoplankton populations. As phytoplankton form the base of the marine food web, this decline can significantly impact marine biodiversity, leading to decreased population sizes in marine species and potential disruptions in ecological balance."

Traditional RAG:

Unstructured information retrieval may lead to missing...



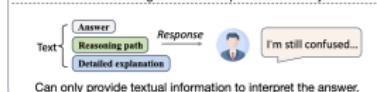
Hard to reason about distributed domain knowledge in long inputs.



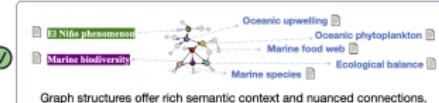
Computationally and time consuming for large vector databases.



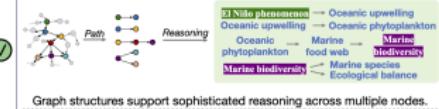
Difficult to integrate with multiple data modality.



Can only provide textual information to interpret the answer.

GraphRAG:

Graph structures offer rich semantic context and nuanced connections.



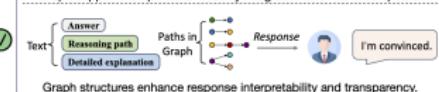
Graph structures support sophisticated reasoning across multiple nodes.



Graph databases optimize for relationship-based queries and faster retrieval.

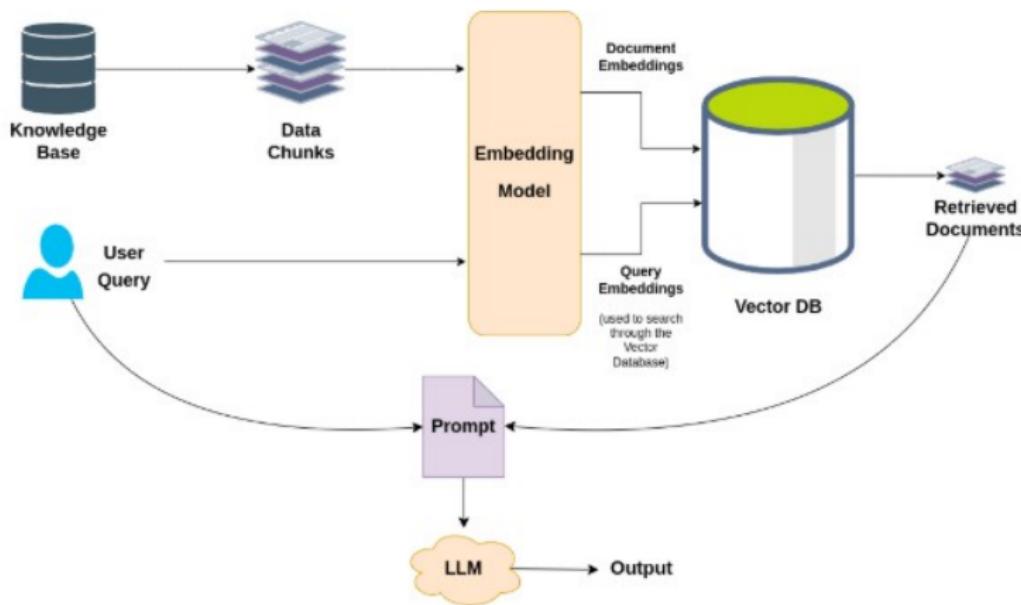


Graph support multiple data modality integration and real-time update.



Graph structures enhance response interpretability and transparency.

Limitations of Traditional RAG



(Ref: GraphRAG: The Practical Guide for Cost-Effective Document Analysis with Knowledge Graphs -Jaykumaran)

Challenges in Standard RAG

- ▶ **Lack of Explainability:** Hard to trace the source of retrieved information.
- ▶ **Local Window:** Limited chunk-level context leads to incomplete responses.
- ▶ **Scalability Issues:** Struggles with large-scale medical and legal corpora.
- ▶ **Loss of Structural Relationships:** Ignores hierarchical relationships in data.

GraphRAG: Evolution of Knowledge Graphs

- ▶ Earlier, KGs were built using statistical and NLP techniques.
- ▶ GraphRAG scales efficiently by using LLMs for entity extraction.
- ▶ Entities in KG are nodes, linked by edges encoding relationships.

Advantages of Knowledge Graphs

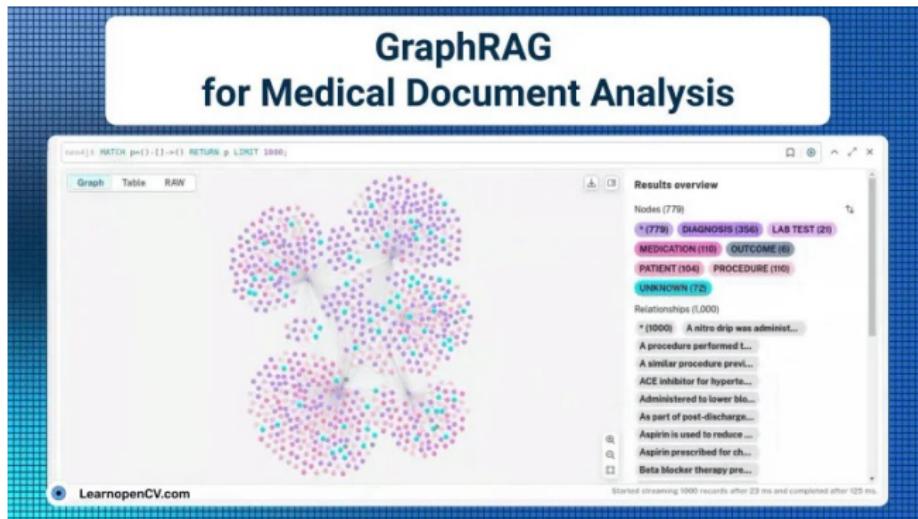
- ▶ Connects information from multiple sources for deeper insights.
- ▶ Boosts retrieval accuracy and enables multi-hop reasoning.
- ▶ Enhances LLM responses by integrating structured relationships.

Key Features of GraphRAG

- ▶ Maintains hierarchical communities preserving local and global insights.
- ▶ Ensures source traceability down to the node level for citations.
- ▶ Aggregates information from multiple sources, reducing bias.
- ▶ Focuses on meaningful nodes, filtering out irrelevant information.

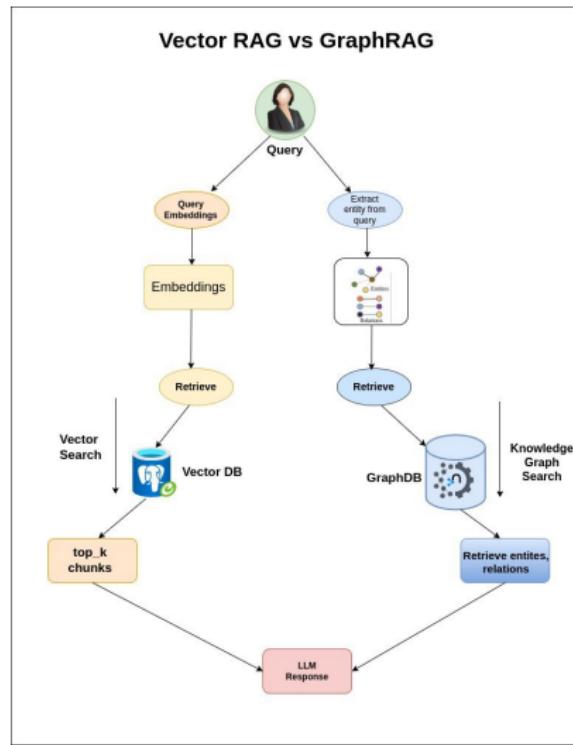
Example: Medical Query Challenge

- ▶ Query: How does Medication A in Patient Record 1 affect Condition B in Patient Record 2?
- ▶ LLM needs to infer relationships across multiple records.
- ▶ Standard RAG struggles with such complex dependencies.
- ▶ Scaling this to millions of patient records is infeasible.



(Ref: GraphRAG: The Practical Guide for Cost-Effective Document Analysis with Knowledge Graphs -Jaykumaran)

Vector RAG vs GraphRAG



(Ref: GraphRAG: The Practical Guide for Cost-Effective Document Analysis with Knowledge Graphs -Jaykumaran)

Traditional RAG vs GraphRAG

Feature	Traditional RAG	GraphRAG
Data Representation	Flat Vectors	Knowledge Graph
Query Scope	Local context	Global Reasoning
Scalability	Low	High
Citation Transparency	Low	High (Traceable sources)
Response Coherence	Fragmented	Relevant and Context-Rich

How GraphRAG Solves the Problem

- ▶ GraphRAG combines graph structures with vector search.
- ▶ Traverses multi-hop connections to infer relationships.
- ▶ Offers more reliable responses than vector-only RAG.

Graph-based RAG helps agentic AI make human-like decisions. – May Habib,
CEO of Writer.com



Applications of GraphRAG

- ▶ **Healthcare:** Assists in diagnoses and treatment decisions.
- ▶ **Banking:** Detects fraudulent transactions using knowledge graphs.
- ▶ **Customer Service :** quickly answer customer questions from thousands of pages of policy documentation
- ▶ **Recommendations:** understand customer behavior and preferences better, to provide personalized services.
- ▶ **Supply Chain:** product recall and associated quality control checking, internal documentation search

How GraphRAG Works?

- ▶ **Knowledge Graph Construction:** Extracts entities and relationships.
- ▶ **Knowledge Graph Summarization:** Generates hierarchical summaries.
- ▶ **Retrieval-Augmented Generation:** Uses local and global searches for queries.

Example of GraphRAG Representation

```
1 # Entities
2 Type 2 Diabetes (Condition)
3 High Blood Sugar Levels (Symptom)
4 Nerve Damage (Complication)
5 Kidney Disease (Complication)
6 Cardiovascular Problems (Complication)
7
8 # Relationships
9 Type 2 Diabetes -> has_symptom -> High Blood Sugar Levels
10 Type 2 Diabetes -> can_lead_to -> Nerve Damage
11 Type 2 Diabetes -> can_lead_to -> Kidney Disease
12 Type 2 Diabetes -> can_lead_to -> Cardiovascular Problems
13
```



Advantages of GraphRAG

- ▶ Structured knowledge representation.
- ▶ Context-aware and efficient retrieval.
- ▶ Handles complex queries effectively.
- ▶ Provides explainability and transparency.

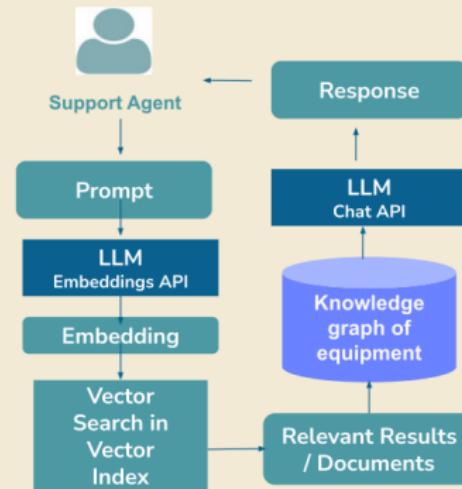
Improve Answers: Pure text w/ Knowledge Graph

Global Equipment Manufacturer:

Challenge: Support agents want to lower their mean time to repair (MTTR) by quickly providing a resolution, but the best answer is hidden among thousands of field docs, engineering summaries, documentation, and case histories.

Solution: Use vector embeddings to look for the most relevant information, while the Neo4j knowledge graph adds the exact match for critical components.

Desired Outcome: Reduce MTTR, and provide the best option for the customer quickly and with confidence.



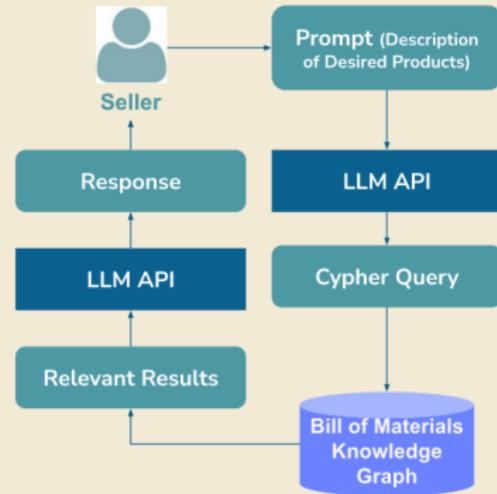
Knowledge Access: Pure data w/ Text to Query

Global Electronics Manufacturer:

Challenge: Sellers want to quote the right product combination and price for every opportunity across five unique types of external stakeholders. Information is stored across millions of bills of materials and product combinations.

Solution: Bring bills of materials into a knowledge graph that can be queried by an LLM interface by sellers.

Desired Outcome: Quote the best option for each customer quickly and with confidence.



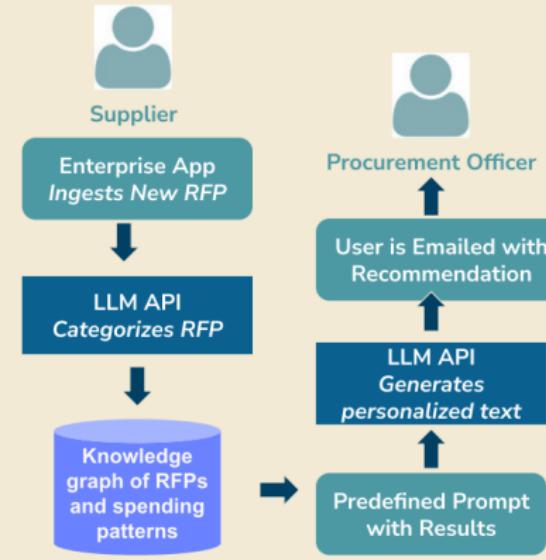
Complex Processes: Mixed Text + Data

Government Procurement Entity:

Challenge: Government entities can overspend or acquire goods and services redundantly because the volume of RFPs makes reviewing each of them resource prohibitive.

Solution: Use LLM to read the nature of RFP and classify it accordingly. Compare the new RFP against the knowledge graph of active and historic RFPs and associated spend. Recommend opportunities for consolidation or terms negotiation with suppliers.

Desired Outcome: Save tax dollars for other important projects

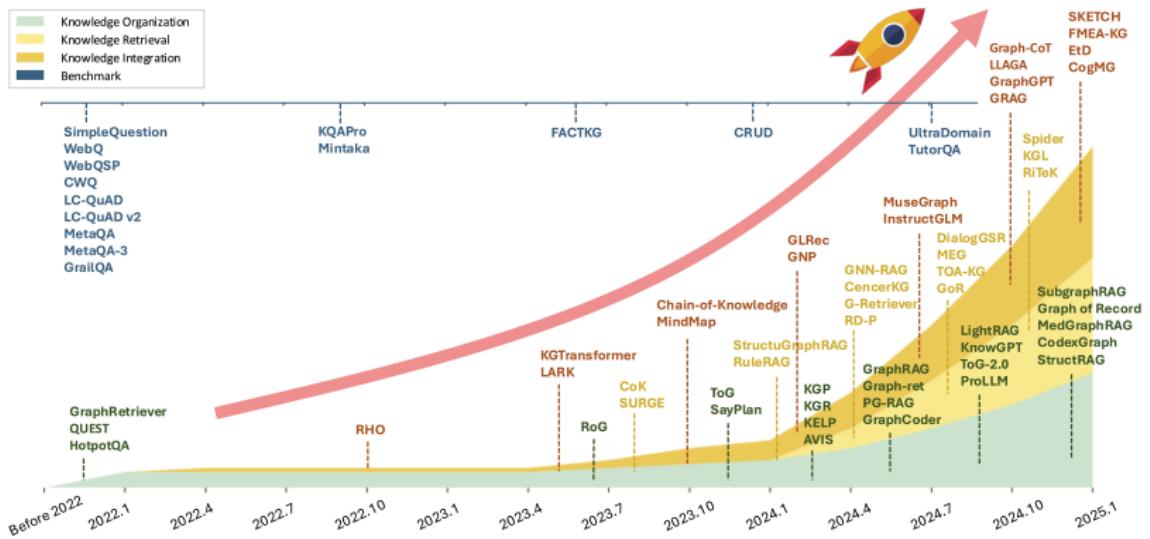


Challenges of GraphRAG

- ▶ **Complex Knowledge Graph Construction:** Requires sophisticated NLP techniques.
- ▶ **Data Dependency:** Performance relies on input data quality.
- ▶ **Scalability Issues:** Large graphs require significant computational resources.



Trend of GraphRAG Research



(Ref: Awesome-GraphRAG (GraphRAG Survey))

Conclusion

- ▶ GraphRAG enhances traditional RAG models using structured knowledge.
- ▶ Improves accuracy, context-awareness, and efficiency.
- ▶ Useful in various domains like healthcare and banking.
- ▶ A promising approach for future AI-powered knowledge retrieval.



Concepts

YHK

Cost-Efficient GraphRAG Approach

- ▶ Leverages graphs for RAG without high costs.
- ▶ Minimizes reliance on LLMs or uses smaller on-premise models.
- ▶ Structured as a layered graph:
 - ▶ Ontology Layer – Defines domain structure (fixed or nearly fixed).
 - ▶ Document Layer – Contains chunked documents, similar to a vector DB.
 - ▶ Entity Layer (Optional) – Extracted entities enhance search.

Challenges in Ontology-Based Graphs

- ▶ Not all datasets belong to a well-defined domain.
- ▶ Subject Matter Experts (SMEs) may not be available.
- ▶ Eliminating fixed ontology layer could improve flexibility.

NLP-Powered Graph Approach

- ▶ Drops the ontology layer for cost-efficiency.
- ▶ Graph consists of:
 - ▶ Document Layer – Contains document chunks.
 - ▶ Tokens Layer – Extracted tokens improve search.
- ▶ NLP reduces dependency on LLMs.



Data Preprocessing Pipeline

- ▶ Chunking – Splitting documents into segments.
- ▶ Embedding – Using Hugging Face model for embeddings.
- ▶ Graph Construction – Built using NetworkX or Neo4j.
- ▶ Token Extraction – Generates token, bigram, and trigram nodes.

Indexing and Interconnectivity

- ▶ Tokens are shared across documents, interlinking content.
- ▶ Need to connect entities using context, logic, and semantics.
- ▶ Avoid reliance on massive models due to cost constraints.

Triple Extraction for Graph Optimization

- ▶ Triple extraction improves retrieval efficiency.
- ▶ Uses a smaller transformer model fine-tuned for this task.
- ▶ Triplets mapped to token nodes enhance query performance.
- ▶ Queries traverse graph using triplet relationships.

Optimizing Graph-Based Retrieval

- ▶ Combines standard RAG with triplet-enhanced retrieval.
- ▶ Retrieves relevant text chunks and connected triplets.
- ▶ Improves search relevance and reduces retrieval complexity.

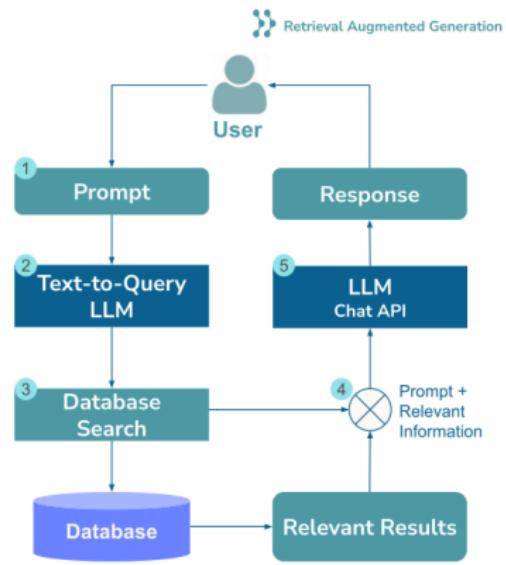
Approaches

YHK

Pure Data Retrieval

1. accept prompt from user
2. pass prompt to a fine-tuned code generation model
3. run query against database
4. combine user prompt with query results
5. generate natural language results with LLM

26 Neo4j Inc. All rights reserved 2023



neo4j

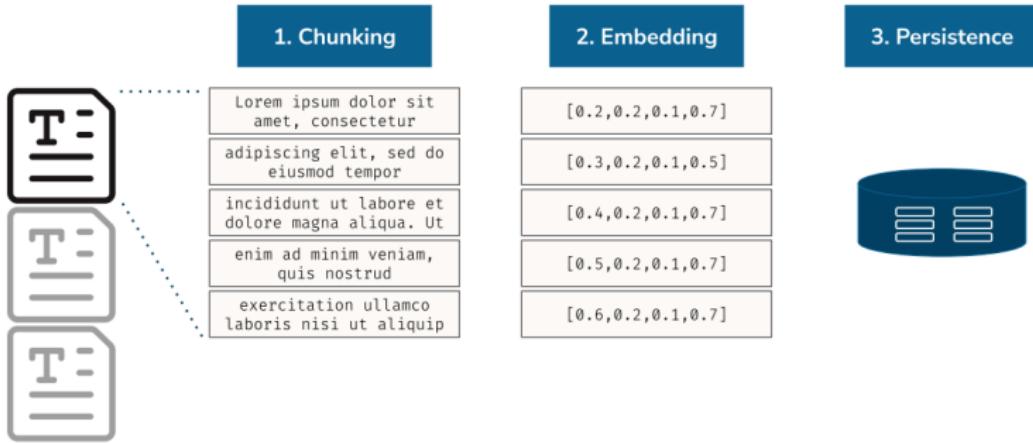
Challenges

- ▶ Getting it to work at all – generating syntactically correct queries
- ▶ Getting it to do the right thing – producing meaningful results
- ▶ Avoiding accidents – mistaken deletion
- ▶ Preventing malicious intent – SQL injection gone wild

(Ref: The GenAI Stack - Andreas Kollegger - Neo4j)



Pure Text Preparation



30 Neo4j Inc. All rights reserved 2023

Challenges

- ▶ How?
 - ▶ pick a chunk method & size
 - ▶ each chunk is a record
 - ▶ store chunk with metadata
 - ▶ connect each chunk to original document
 - ▶ connect previous/next chunk
- ▶ Challenges:
 - ▶ what makes a good chunk?
 - ▶ potential chunk duplication
 - ▶ how to re-assemble chunk context?
 - ▶ what about cross-document chunks?
 - ▶ explaining the relevance

(Ref: The GenAI Stack - Andreas Kollegger - Neo4j)



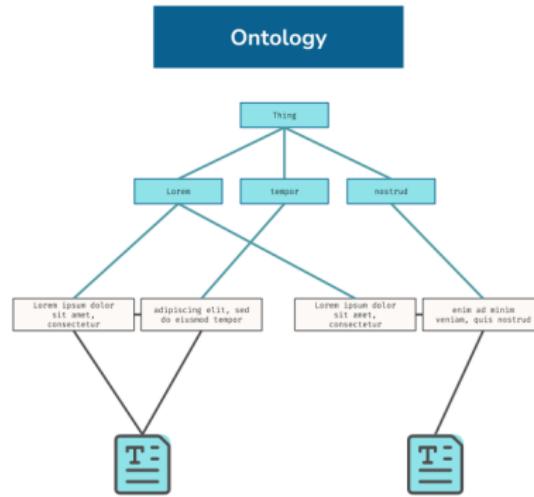
Explicit Similarity

How?

- named entity recognition
- metadata extraction
- document cross-linking, page ranking

Challenges:

- does chunk similarity provide the most relevant answer?
- what about the person asking the question?



Mixed Text & Data

How?

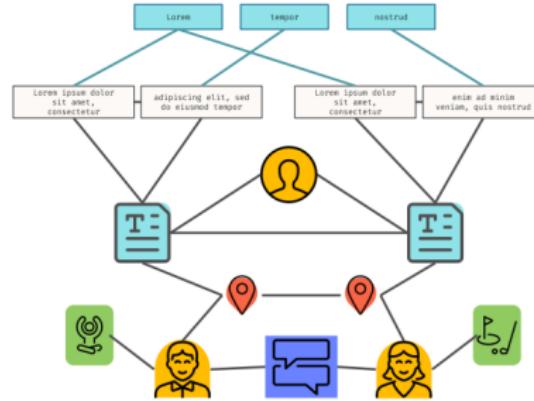
- extracted information
- application data
- user data
- explicit semantic connections

Challenges:

- what is most relevant, to the user?

39 Neo4j Inc. All rights reserved 2023

Mixed Text & Data



Knowledge Graphs

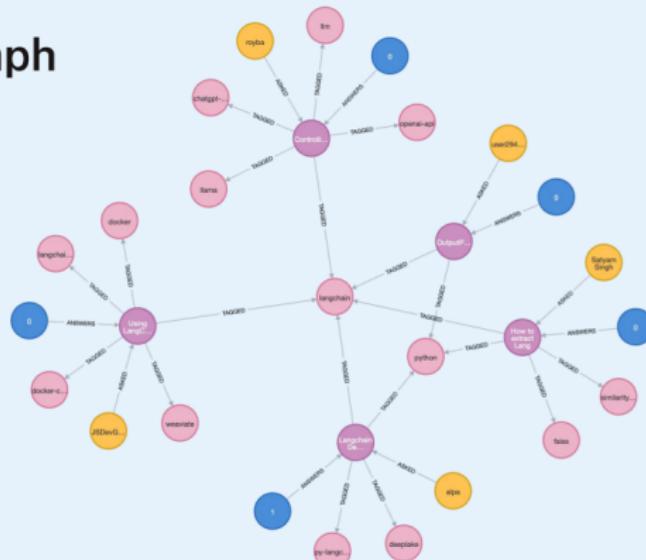
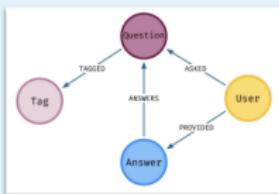
YHK

DEFINITION

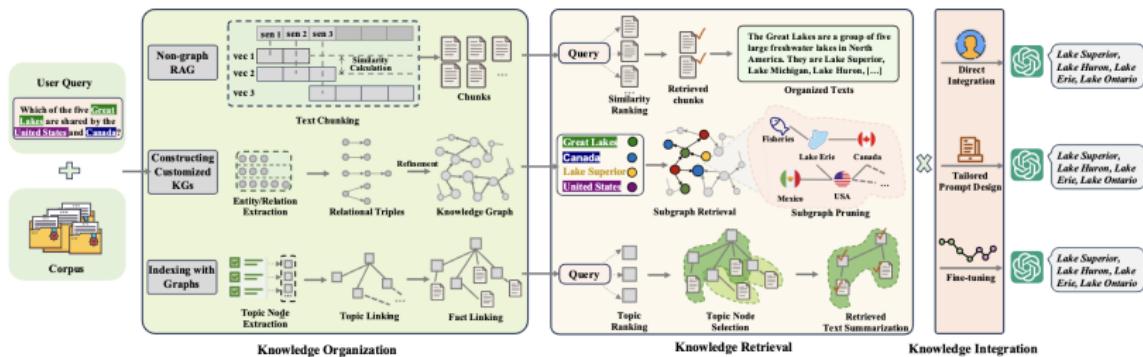
A Knowledge Graph is a structured way of representing information, typically using nodes and edges to depict relationships between entities (e.g., people, places, things, concepts). These entities and their interconnections form a graph-like structure, which can be used to model complex sets of data and the relationships within that data.

(Ref: The GenAI Stack - Andreas Kollegger - Neo4j)

StackOverflow Knowledge Graph



49 Neo4j Inc. All rights reserved 2023



Overview

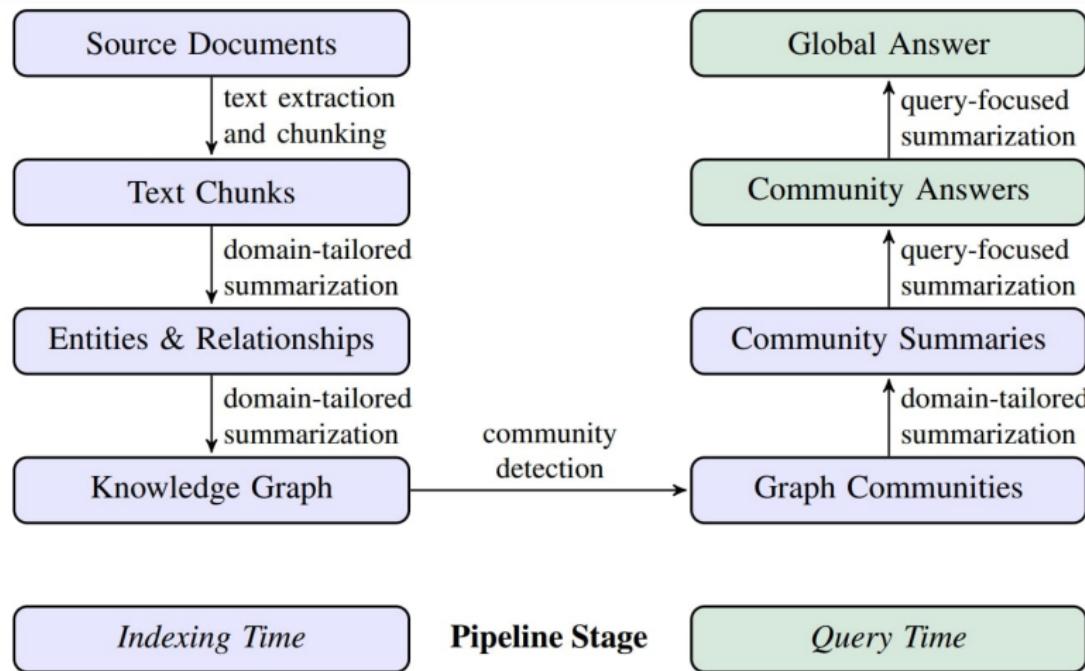
Overview of traditional RAG and two typical GraphRAG workflows.

- ▶ Non-graph RAG organizes the corpus into chunks, ranks them by similarity, and retrieves the most relevant text for generating responses.
- ▶ Knowledge-based GraphRAG extracts detailed knowledge graphs from the corpus using entity recognition and relation extraction, offering fine-grained, domain-specific information.
- ▶ Index-based GraphRAG summarizes the corpus into high-level topic nodes, which are linked to form an index graph, while the fact linking maps topics to text.

(Ref: Awesome-GraphRAG (GraphRAG Survey))



Microsoft GraphRAG Workflow



(Ref: GraphRAG: The Practical Guide for Cost-Effective Document Analysis with Knowledge Graphs -Jaykumaran)

Indexing Phase Overview

- ▶ Two-stage process: Generate Knowledge Graph (KG) and Community Clustering.
- ▶ Extract entities and relationships, construct KG.
- ▶ Partition KG into semantically related clusters.
- ▶ Generate community summaries for optimized retrieval.

Step 1: Segment Documents into Chunks

- ▶ Large text split into smaller chunks for LLM processing.
- ▶ Small chunks retain details but increase API cost.
- ▶ Larger chunks reduce cost but may omit key data.
- ▶ Self-reflection in prompts improves extraction accuracy.



Step 2: Extract Entities and Relationships

- ▶ LLM extracts entities, relationships, and descriptions.
- ▶ Entities assigned unique IDs for traceability.
- ▶ Pronouns and ambiguities resolved.

```
1  {
2      "nodes": [
3          {"id": "n0", "entity_type": "PATIENT", "name": "71-YEAR-OLD
4              GENTLEMEN"},
5              {"id": "n1", "entity_type": "DIAGNOSIS", "name": "PERIPHERAL VASCULAR
6                  DISEASE"}
7          ],
8          "edges": [
9              {"id": "e0", "source": "n0", "target": "n1", "description": "Patient
10                 has history of this condition"}
11         ]
12     }
```

Step 2: Extract Entities and Relationships

```
1  **-Goal-**  
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 identified entities.  
3  
-Steps-  
5  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: [{entity_types}]  
    - entity_description: Comprehensive description of the entity's  
9  
2. From the entities identified in step 1, identify all pairs of  
    (source_entity, target_entity) that are *clearly related to each other.  
11  
For each pair of related entities, extract the following information:  
13  - source_entity: name of the source entity, as identified in step 1  
    - target_entity: name of the target entity, as identified in step 1  
15  - 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  
17
```

Step 2: Extract Entities and Relationships

```
# Entities:  
2 {  
    "nodes": [  
        {  
            "id": "n0",  
            "entity_type": "PATIENT",  
            "description": "A male patient aged 71 with several medical conditions",  
            "source_id": "n0",  
            "name": "71-YEAR-OLD GENTLEMEN"  
        },  
        {  
            "id": "n1",  
            "entity_type": "DIAGNOSIS",  
            "description": "Condition treated with right leg angioplasty\nPeripheral vascular disease (PVD) is a circulatory condition characterized by narrowed or blocked arteries outside of the heart and brain, leading to reduced blood flow, particularly to limbs. Patients with PVD may have undergone procedures like angioplasty in affected areas to improve circulation. Treatment generally involves medications such as Plavix and aspirin, as well as surgical interventions based on individual assessments.",  
            "source_id": "n1",  
            "name": "PERIPHERAL VASCULAR DISEASE"  
        },  
    ]  
}
```



Step 2: Extract Entities and Relationships

```
1 # Relations
2   "edges": [
3     {
4       "id": e0,
5       "source": "n0",
6       "target": "n1",
7       "description": "The patient has a medical history of this condition",
8     },
9     {
10       "id": e1,
11       "source": "n0",
12       "target": "n2",
13       "description": "Chronic Obstructive Pulmonary Disease is part of his
14         health profile",
15     },
16   ]
```

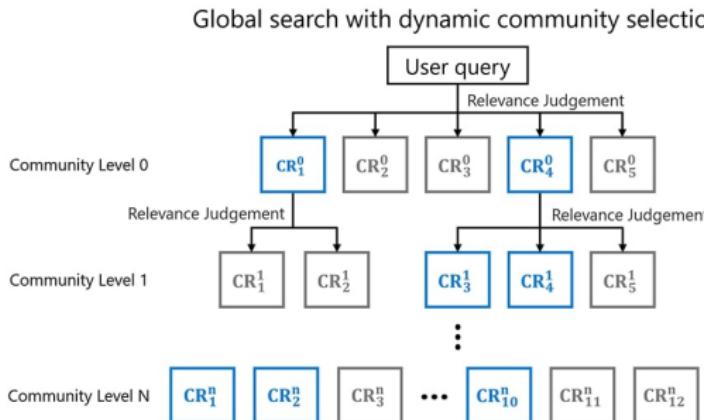


Step 3: Construct Knowledge Graph

- ▶ Entities and relationships aggregated into a unified KG.
- ▶ Duplicates merged for efficiency.

Step 4: Graph Partitioning into Communities

- ▶ Leiden algorithm partitions KG into clusters.
- ▶ Hierarchical approach ensures fine granularity.
- ▶ Nodes belong to only one community.
- ▶ Dynamic community selection at multiple levels.



Step 1: Recursive Community Selection

Start at root level (level 0), using an LLM to rate the relevance of each community report (CR) against the user query.

If the community report is not relevant (grey boxes), then we move on to the next report.

If the community report is relevant, then we repeat the same operation to its child nodes. We keep the report (blue boxes) if none of its child communities is relevant or if it is a leaf node.

Step 2: Map-reduce Search

Perform map-reduce search on the selected community reports to generate the final response.

(Ref: GraphRAG: The Practical Guide for Cost-Effective Document Analysis with Knowledge Graphs -Jaykumaran)

Step 5: Community Summarization

```
# **Goal**  
2 Write a comprehensive report of a community, given a list of  
entities that belong to the community as well as their  
relationships and optional associated claims. The report  
will be used to inform decision-makers about information  
associated with the community and their potential impact.  
The content of this report includes an overview of the  
community's key entities, their legal compliance,  
technical capabilities, reputation, and noteworthy claims.
```

4

YHK

Step 5: Community Summarization

```
1 # Report Structure
```

The report should include the following sections:

- ```
3 - TITLE: community's name that represents its key entities
```
- ```
5 - SUMMARY: An executive summary of the community's overall
```
- ```
 structure,
```
- ```
7 - IMPACT SEVERITY RATING: a float score between 0-10 that
```
- ```
 represents the severity of IMPACT posed by entities
```
- ```
    within the community.
```
- ```
9 - RATING EXPLANATION: Give a single sentence explanation of
```
- ```
    the IMPACT severity rating.
```
- ```
- DETAILED FINDINGS: A list of 5-10 key insights about the
```
- ```
    community.
```

Querying Phase Overview

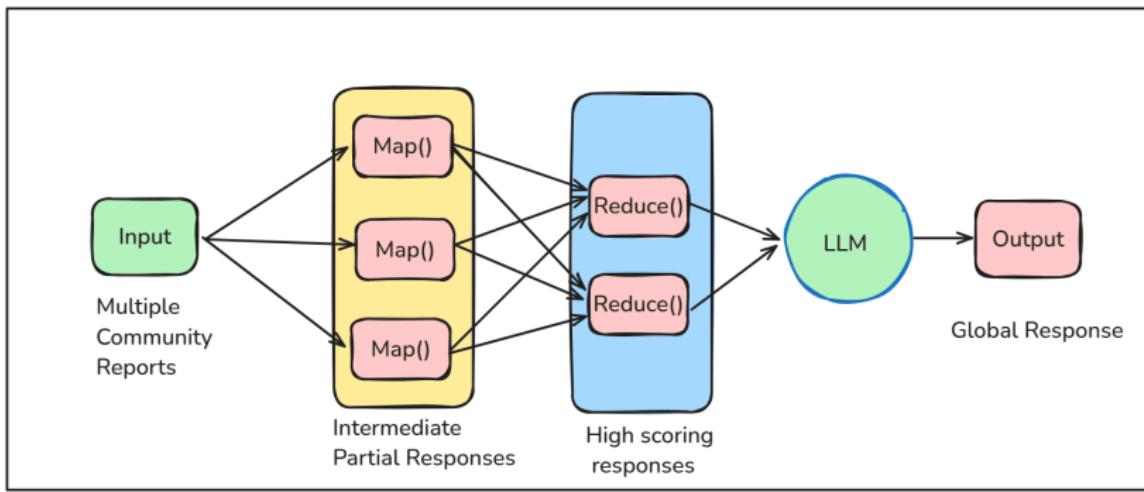
- ▶ Identifies key entities and relationships in queries.
- ▶ Retrieves relevant community summaries instead of entire KG.

Query Focused Summarization (QFS)

- ▶ Identifies patterns and correlations from queries.
- ▶ Matches against community summaries.
- ▶ Optimized retrieval ensures efficient responses.

Map Reduce Approach

- ▶ **Map Phase:** Generates multiple partial responses with importance scores.
- ▶ **Reduce Phase:** Sorts and merges top-scoring responses into a global answer.
- ▶ Efficiently filters search space while maintaining accuracy.



(Ref: GraphRAG: The Practical Guide for Cost-Effective Document Analysis with Knowledge Graphs -Jaykumaran)

YHK

Benchmarks and Evaluation Metrics

- ▶ Traditional RAGAS metrics are inadequate for sensemaking queries.
- ▶ LLM judges responses based on:
 - ▶ **Comprehensiveness:** Covers all aspects of the query.
 - ▶ **Diversity:** Provides varied perspectives and insights.
 - ▶ **Empowerment:** Aids decision-making without false assumptions.
 - ▶ **Directness:** Clear, concise, and to the point.
- ▶ GraphRAG excels in comprehensiveness and diversity.
- ▶ Naive RAG is better at generating direct answers.

Shortcomings of Microsoft GraphRAG

- ▶ Multiple LLM API calls make it slow, inefficient, and costly.
- ▶ No de-duplication step, resulting in a noisy graph index.
- ▶ Upserting new data requires reconstructing the entire KG, making it impractical for production.

GraphRAG Alternatives and Improvements

- ▶ Graph-based RAG is not new; earlier versions by NebulaGraph, Langchain, LlamaIndex, Neo4j.
- ▶ Microsoft's approach became widely accepted but has limitations.
- ▶ Led to new variants focusing on efficiency and accuracy.



GraphRAG Variants

- ▶ **Nano-GraphRAG:** Introduced lighter, faster versions.
- ▶ **Top-k Selection:** Unlike map-reduce, it selects only top-k communities for efficiency.
- ▶ Derived versions: LightRAG, MedGraphRAG, FastGraphRAG.
- ▶ MedGraphRAG is challenging to adapt to OpenAI endpoints.



FastGraphRAG Performance

- ▶ 27x faster and 40% more accurate retrieval.
- ▶ Benchmarks (2wikimultihopQA, 101 queries):

Method	Accuracy (All)	Accuracy (Multihop)	Insertion Time (min)
VectorDB	0.42	0.23	~0.3
LightRAG	0.45	0.28	~25
GraphRAG	0.73	0.64	~40
FastGraphRAG	0.93	0.90	~1.5

(Ref: GraphRAG: The Practical Guide for Cost-Effective Document Analysis with Knowledge Graphs -Jaykumaran)

Graph Construction Economy Principle

- ▶ Complex graphs require significant computational resources.
- ▶ Trade-off: Performance gains must justify resource investment.
- ▶ Aim: Maximize performance-to-resource ratio.

FastGraphRAG at Inference Time

- ▶ Uses a PageRank-like algorithm (similar to Google).
- ▶ Determines importance of elements in the KG.
- ▶ Retrieves only the most relevant entities and relationships.
- ▶ Produces high-quality, precise responses.

Evaluating FastGraphRAG for Your Use Case

- ▶ Consider cost-to-performance ratio for your specific dataset.
- ▶ Assess viability based on document complexity and enterprise needs.
- ▶ FastGraphRAG offers a scalable approach with efficiency gains.



Conclusion

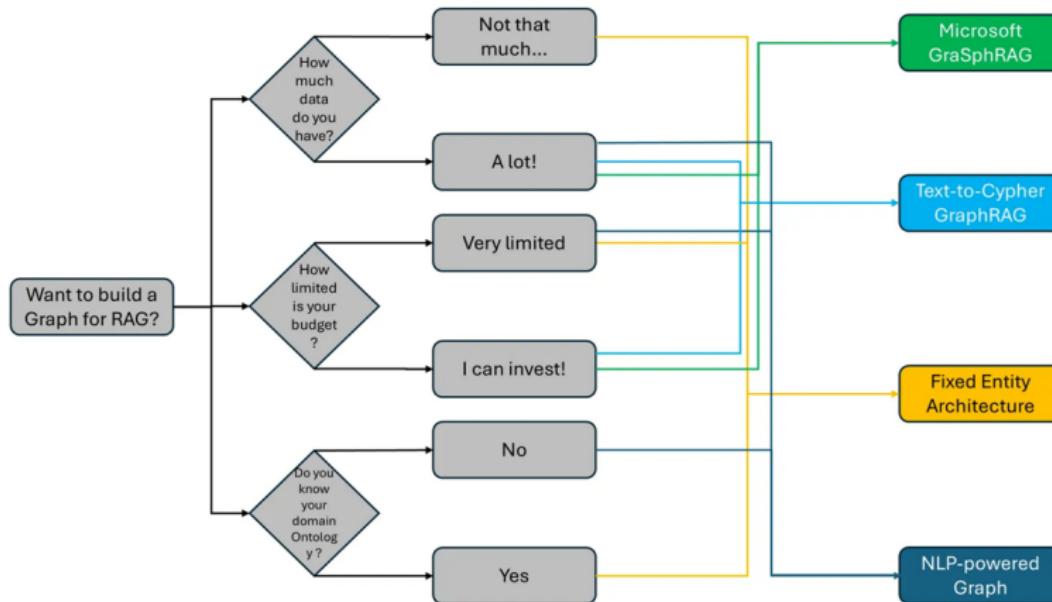
- ▶ GraphRAG scales from local to global knowledge representation.
- ▶ Community-based summaries improve response quality.
- ▶ Enables optimized retrieval with hierarchical partitioning.



Implementation

YHK

How to build GraphRAG?



(Ref: Build your hybrid-Graph for RAG & GraphRAG applications using the power of NL - Irina Adamchic)



Neo4j Ollama GraphRAG



(Ref: The GenAI Stack - Andreas Kollegger - Neo4j)

- ▶ Application: LangChain: an “orchestration framework” for integrating with LLMs
- ▶ Gen AI : Ollama : locally managed language models
- ▶ Database: Neo4j: knowledge graph to augment the LLM
- ▶ Files: import: sample data sources for constructing a knowledge graph

Code Walkthrough: Medical Doc Analysis with Fast GraphRAG

(Ref: GraphRAG: The Practical Guide for Cost-Effective Document Analysis with Knowledge Graphs -Jaykumaran)



FastGraphRAG Overview

- ▶ All functions are asynchronous, enabling concurrent API calls.
- ▶ Works well with providers without RPM limits (e.g., OpenAI).

Installation

```
1 !git clone https://github.com/circlemind-ai/fast-graphrag.git
2 %cd fast_graphrag
3 %poetry install
4
5 # or
6 !pip install fast-graphrag
```

Importing Dependencies

```
2 import instructor # Structured outputs for LLMs
3 import os
4 from fast_graphrag import GraphRAG, QueryParam
5 from fast_graphrag._llm import OpenAIEmbeddingService, OpenAILLMService
```

Dataset: MIMIC-IV-ICU

- ▶ Public dataset of de-identified ICU patient EHRs.
- ▶ Using first 430 samples (650 words per file).

```
2      mimic_ex_500/
3          |- report_0.txt
4          |- report_1.txt
5          |- ...
6          |- report_499.txt
7          |- report_500.txt
```

Defining Domain and Entity Types

```
DOMAIN = "Analyze clinical records for key medical
entities."
EXAMPLE_QUESTIONS = [
    "What are the common risk factors for sepsis in ICU
patients?",  
4
    "How do trends in lab results correlate with patient
outcomes in cases of acute kidney injury?",  
5
    "Describe the sequence of interventions for patients
undergoing major cardiac surgery.",  
6
    "How do patient demographics and comorbidities influence
treatment decisions in the ICU?",  
7
    "What patterns of medication usage are observed among
patients with chronic obstructive pulmonary disease
(COPD)?"]  
8
ENTITY_TYPES = ["Patient", "Diagnosis", "Procedure", "Lab
Test", "Medication", "Outcome", "Unknown"]
working_dir = "./WORKING_DIR/mimic_ex500/"  
10
```

Configuring the FastGraphRAG Pipeline

```
1 grag = GraphRAG(
2     working_dir=working_dir,
3     n_checkpoints=2, # to save backups (recommended)
4     domain=DOMAIN,
5     example_queries="\n".join(EXAMPLE_QUERIES),
6     entity_types=ENTITY_TYPES,
7     config=GraphRAG.Config(
8         llm_service=OpenAILLMService(
9             model="Phi4_6k", # gemini-2.0-flash-exp
10            # or https://generativelanguage.googleapis.com/v1beta/openai/
11            base_url="http://localhost:11434/v1", # Ollama
12            api_key="ollama", # or GEMINI_API_KEY
13            mode=instructor.Mode.JSON,
14            client="openai"),
15         embedding_service=OpenAIEmbeddingService(
16             model="mxbai-embed-large", # mxbai-embed-large
17             base_url="http://localhost:11434/v1",
18             api_key="ollama",
19             embedding_dim=1024, # mxbai-embed-large - 1024 ; nomic-embed - 768
20             client="openai"),
21         ),
22     )
23 directory_path = "mimic_ex_430" # input dir to the dataset
```

Graph Indexing

```
1 def graph_index(directory_path):
2     file_count = 0 # Keep track of processed files.
3     for filename in os.listdir(directory_path):
4         if filename.endswith('.txt'):
5             file_path = os.path.join(directory_path, filename)
6             with open(file_path, 'r') as file:
7                 content = file.read()
8                 if isinstance(content, list):
9                     content = "\n".join(map(str, content))
10                if isinstance(content, dict):
11                    key_to_use = next(iter(content.keys()), None)
12                    content = content[key_to_use] if key_to_use else
13                        str(content)
14                else: content = str(content)
15                graph.insert(content)
16                file_count += 1
17                total_files = sum(1 for f in os.listdir(directory_path) if
18                    f.endswith(".txt"))
19                    print("***** $$$$$ $$$$$ *****")
20                    print(f"Total Files Processed: -> {file_count} / {total_files}")
21                    print("***** $$$$$ $$$$$ *****")
22    return None
23 graph_index(directory_path)
```



Indexing Performance

- ▶ Indexed 430 items in 510 minutes (8.5 hours).
- ▶ Processed locally with a 14B model.

Querying Example 1: Targeted at Global

```
1 grag.query("Describe the sequence of interventions for  
2 cardiac surgery?").response  
  
3 ****  
4 The sequence of interventions for the clinical scenario in  
5 post-aortic valve repair surgery with complications like  
6 respiratory issues, neurological impairment, and renal  
7 failure includes:  
8 **Pre-surgery Preparation:**  
9 - Evaluate and stabilize the patient.  
10 **Intraoperative Procedures:**  
11 - Conduct aortic valve replacement or repair.  
12 **Immediate Postoperative Care:**  
13 - Continue invasive mechanical ventilation, especially in  
14 patients with respiratory issues or neurological  
15 impairments preoperatively.
```

11

YHK

Querying Example 1: Targeted at Global

```
**Rehabilitation and Follow-up:**  
2 - Perform swallow evaluation tests to ensure readiness for  
    oral intake, followed by a diet progression to soft or  
    clear liquids if necessary.  
  
4 **Commonly Used Medications:**  
- Anticoagulants (e.g., heparin) to prevent blood clots.  
6  
A multidisciplinary team including cardiologists,  
pulmonologists, and other specialists is often involved  
to address any comorbid conditions.  
8 *****  
  
10 # o3-mini-high rated this response : 8/10.
```

Querying Example 2: Targeted at Global

```
grag.query("How do demographics and comorbidities  
influence ICU decisions?").response  
  
*****  
Patient demographics and comorbidities significantly  
influence treatment decisions in the ICU, as evidenced by  
detailed hospital courses described in various scenarios.  
  
1. **Demographics**: Specific ages or life stages, such as  
"on day of life 45" for a patient with unique medical  
conditions (metabolic issues), indicate tailored care  
plans sensitive to developmental stages.  
  
2. **Comorbidities**:  
- A history like short gut syndrome and colectomy  
influences decisions about when and how to restart total  
parenteral nutrition (TPN). \ldots .
```

10

YHK

Querying Example 2: Targeted at Global

- 1 3. ****Interactions**:** Decisions like replacing a PICC line with a Hickman catheter after infection shows proactive steps in avoiding repeat complications due to prior infections from central lines.
4. ****Treatment Customization**:** The management decisions, such as adjusting insulin administration for endocrine challenges even when insulin levels are normal, demonstrate personalized ICU therapy adapting for nuanced metabolic control in conditions like hypoglycemia of unknown origin.

3 These scenarios highlight that ICU treatment is deeply influenced by a patient's demographic specifics \ldots and chronic health issues effectively.

5 *****

o3-mini-high rated this response : 8/10.

7

YHK

Querying Example 3: Targeted at Global

```
1 grag.query("How can lung function be improved in COPD +  
2 Heart Failure?").response  
*****  
3 To optimize the treatment for patients with both COPD and  
4 heart failure, it is crucial to manage fluids carefully  
using diuretics without causing respiratory issues. This  
necessitates precise adjustments in medication dosages,  
closely monitoring fluid balance, and ensuring proper  
oxygenation levels. Right heart catheterization can  
assist in measuring cardiac function and tailoring  
treatments for these patients. The goal is to alleviate  
symptoms of both conditions while preventing  
exacerbations, which often involves using supplemental  
oxygen or bronchodilators as required.  
*****  
6 # o3-mini-high rated this response : 8/10.
```

Querying Example 4: Targeted at Local

```
1 grag.query("Discuss Hypocapnic and Hypoxemic Respiratory  
2 Failure").response  
  
3 Hypocapnic Respiratory Failure and Hypoxemic Respiratory  
4 Failure are two distinct categories of respiratory  
failure, each characterized by different clinical  
features. Conversely,\ldots Accurate assessment and  
differentiation between these conditions are crucial for  
appropriate management and treatment.  
5  
6 In medical practice, both types may coexist or require  
7 different clinical interventions based on the specific  
8 underlying pathology.  
*****  
# o3-mini-high rated this response : 6/10.
```

Limitations of FastGraphRAG

- ▶ Asynchronous and Pydantic-validated, challenging for beginners.
- ▶ 7B-class models struggle with structured responses.
- ▶ Limited documentation and community support.
- ▶ LightRAG offers better support for local LLMs.
- ▶ Faster than GraphRAG/LightRAG but slower than vector-only RAG.
- ▶ Best for large datasets and high-stakes domains.
- ▶ For simple use cases, vector-only RAG is recommended.

Trivia: PDF Indexing Challenges

- ▶ Faced issues with PDF indexing using pdf-plumber.
- ▶ Pydantic validation errors encountered.
- ▶ Works fine with .txt files converted from PDFs.
- ▶ Did not extensively debug the issue.



Conclusion

- ▶ GraphRAG enables tracking patient behavioral patterns.
- ▶ Supports complex, multi-hop queries with accurate responses.
- ▶ Ideal for high-stakes domains requiring precision.
- ▶ Knowledge graphs improve reasoning and interpretability.
- ▶ Hybrid approaches (KG + vectors) outperform naive RAG.

References

Slides primarily borrowed from ...

- ▶ End-to-End Implementation of GraphRAG(Knowledge Graphs + Retrieval Augmented Generation): Unlocking LLM Discovery on Narrative Private Data - Vinod Kumar G R
- ▶ GraphRAG tutorials by Data Science in your pocket
- ▶ GraphRAG: Enhancing LLMs with Knowledge Graphs - Vijayakumar Ramdoss
- ▶ Awesome-GraphRAG (GraphRAG Survey)
<https://github.com/DEEP-PolyU/Awesome-GraphRAG>
- ▶ Build your hybrid-Graph for RAG & GraphRAG applications using the power of NL - Irina Adamchic
- ▶ The GenAI Stack - Andreas Kollegger - Neo4j
- ▶ GraphRAG: The Practical Guide for Cost-Effective Document Analysis with Knowledge Graphs JaykumaranJaykumaran

Thanks ...

- ▶ Search "**Yogesh Haribhau Kulkarni**" on Google and follow me on LinkedIn and Medium
- ▶ Office Hours: Saturdays, 2 to 5pm (IST); Free-Open to all; email for appointment.
- ▶ Email: yogeshkulkarni at yahoo dot com

(<https://www.linkedin.com/in/yogeshkulkarni/>, QR by Hugging Face

QR-code-AI-art-generator, with prompt as "Follow me")

