# Incorporating KGs to Enhance LLM Performance

Máté Vass

Intelligent Human-Computer Interactions

2025-04-09

# Knowledge graphs

# Drawbacks of LLMs

- Not explainable

- Lack of real understanding

- No domain knowledge

- Not up-to-date

- Hallucinations

# Retrieval Augmented Generation

- Add knowledge to LLM context

  ○ Context window is a limit

- Model editing

  ○ Computanional cost
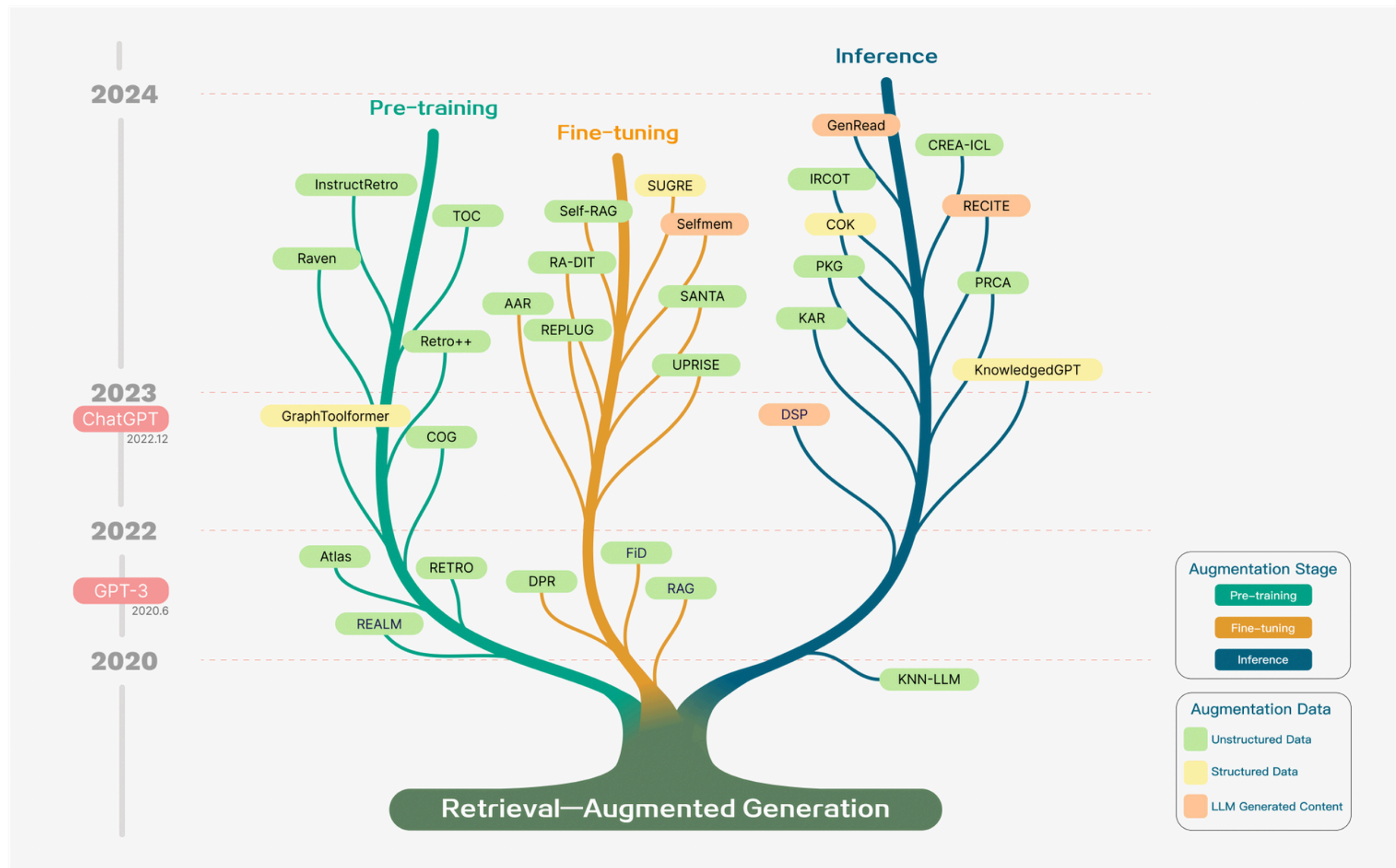
- Extracting relevant knowledge is crucial!

Figure 1: A timeline of existing RAG research. The timeline was established mainly according to the release date.

# Naive RAG Pipeline

- Data indexing (text)

- Chunking (split)

- Embedding (chunk - vector pairs)

- Retrieve relevant chunks

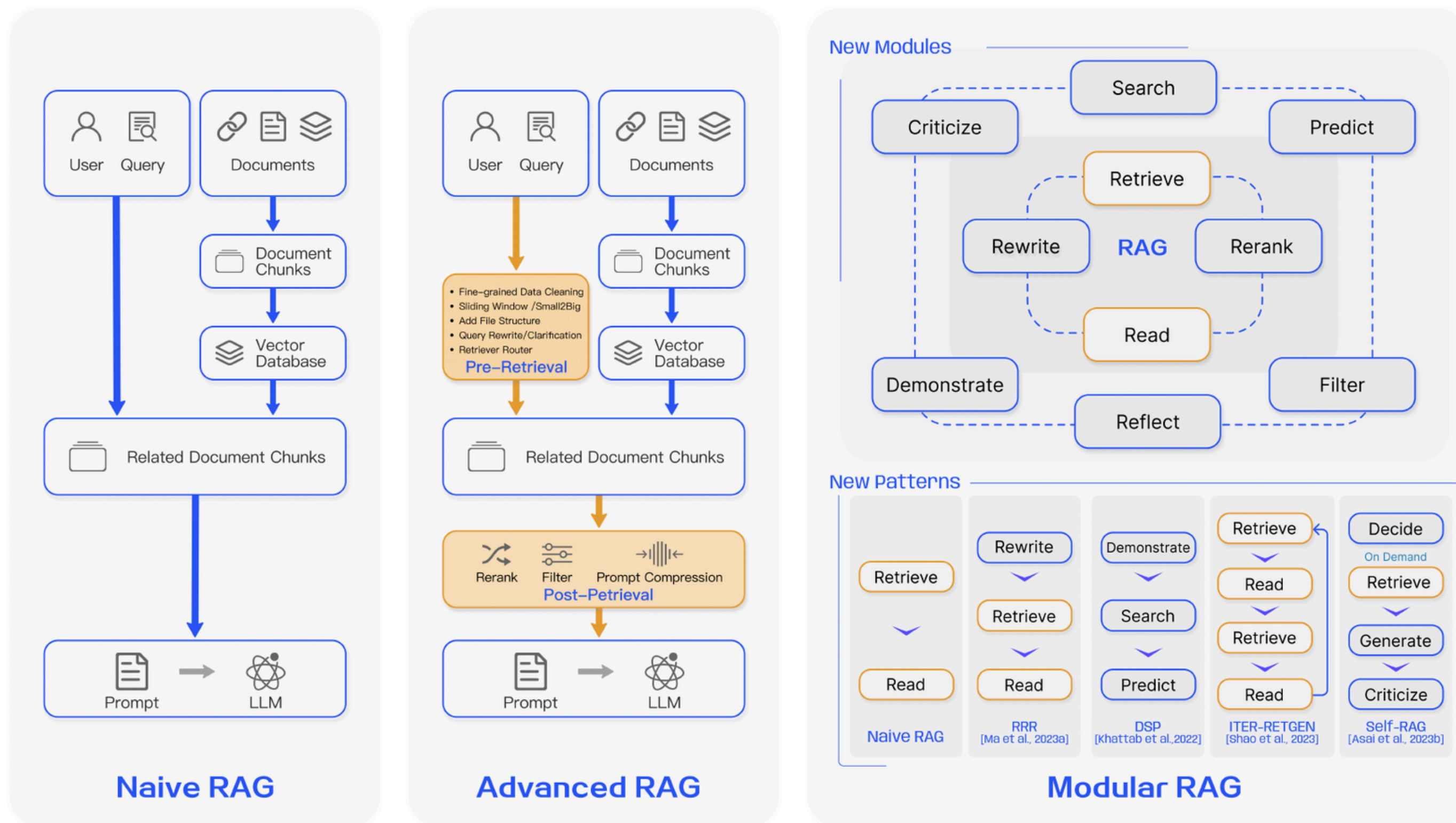- Generate answer using extracted data

Figure 3: Comparison between the three paradigms of RAG

# Examples

- KG is already available

  - KG-GPT

  - KELP
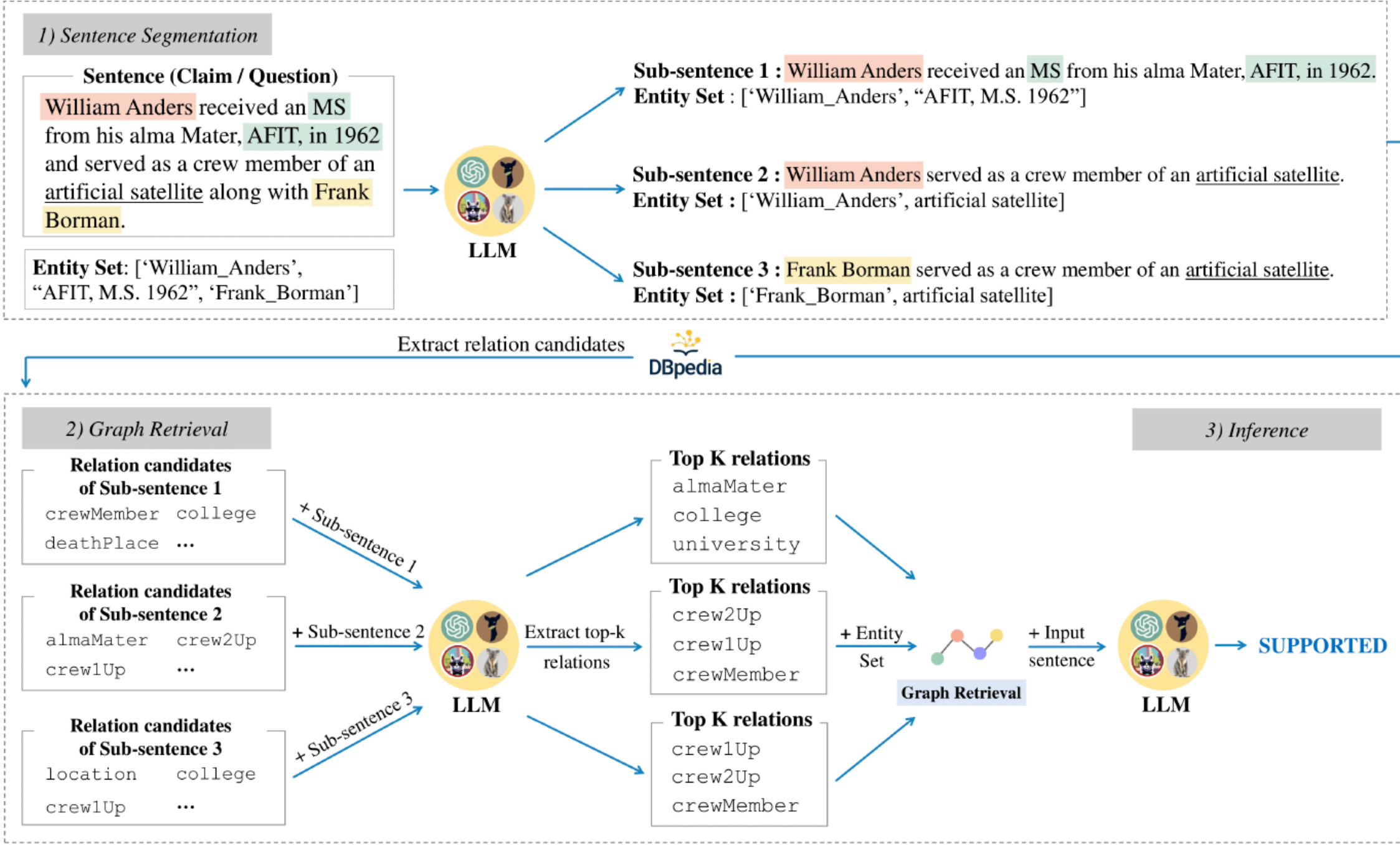
- Textual documents

  - GraphRAG

# KG-GPT



Figure 1: An overview of KG-GPT. The framework comprises three distinct phases: Sentence Segmentation, Graph Retrieval, and Inference. The given example comes from FACTKG. It involves a 2-hop inference from 'William_Anders' to 'Frank_Borman', requiring verification through an evidence graph consisting of three triples. Both 'William_Anders' and 'Frank_Borman' serve as internal nodes in DBpedia (Lehmann et al., 2015), while "AFIT, M.S. 1962" acts as a leaf node. Moreover, *artificial satellite* represents the *Type* information absent from the provided entity set.
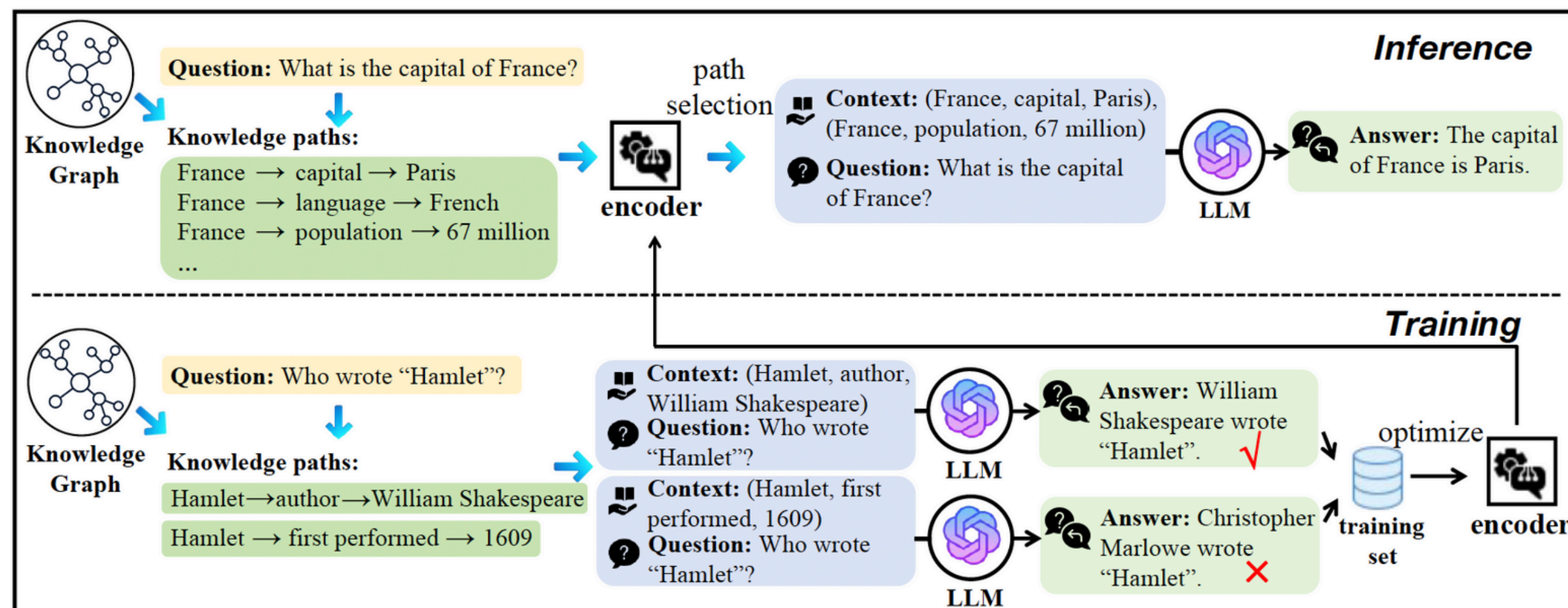
# KELP



Figure 2: The overall pipeline of the proposed KELP. During the inference phase, we identify knowledge paths from the knowledge graph that are associated with the entities present in the input question. An encoder is then trained to select valuable paths as knowledge contexts. Finally, the selected knowledge contexts, along with the input question, are input into the LLM to generate the final answer.

# GraphRAG

- Extracting a knowledge graph out of raw text

- Building a community hierarchy

- Generating summaries for these communities

- Using these structures for RAG-based tasks

# GraphRAG Indexing

- Split text into chunks

- Extract entities, relations

- Hierarchical Leiden community detection

- Bottom-up summary generation
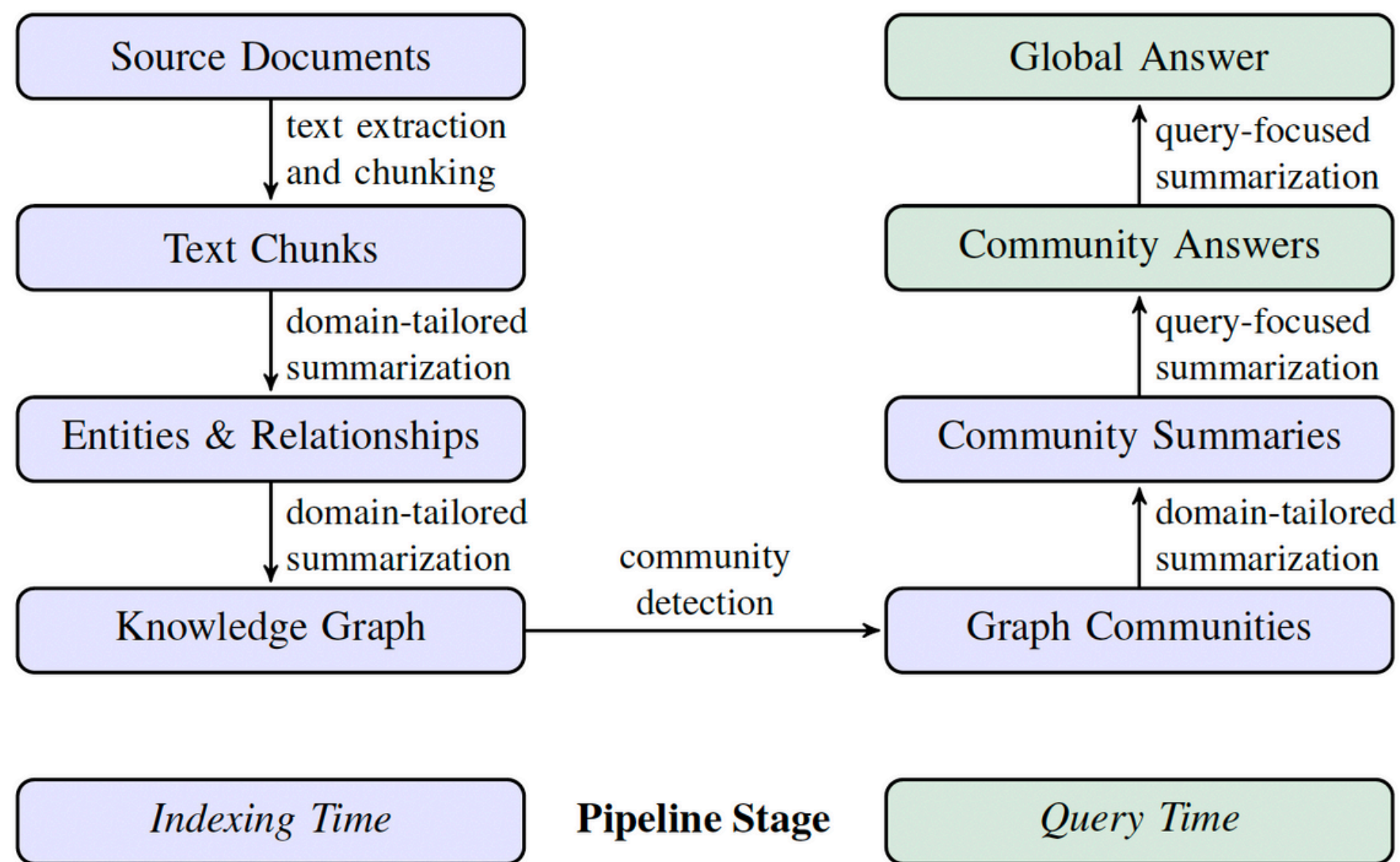
- Query: utilizes the resulting KG

Figure 1: Graph RAG pipeline using an LLM-derived graph index of source document text. This graph index spans nodes (e.g., entities), edges (e.g., relationships), and covariates (e.g., claims) that have been detected, extracted, and summarized by LLM prompts tailored to the domain of the dataset. Community detection (e.g., Leiden, Traag et al., 2019) is used to partition the graph index into groups of elements (nodes, edges, covariates) that the LLM can summarize in parallel at both indexing time and query time. The "global answer" to a given query is produced using a final round of query-focused summarization over all community summaries reporting relevance to that query.

# Thank you for your attention!

Questions?