

CONTENTS

1. RAG Overview
2. RAG Paradigms Shifting
3. Key Technologies and Evaluation
4. RAG Stack and Industry Practices
5. Summary and Prospect

PART 01

Overview of RAG

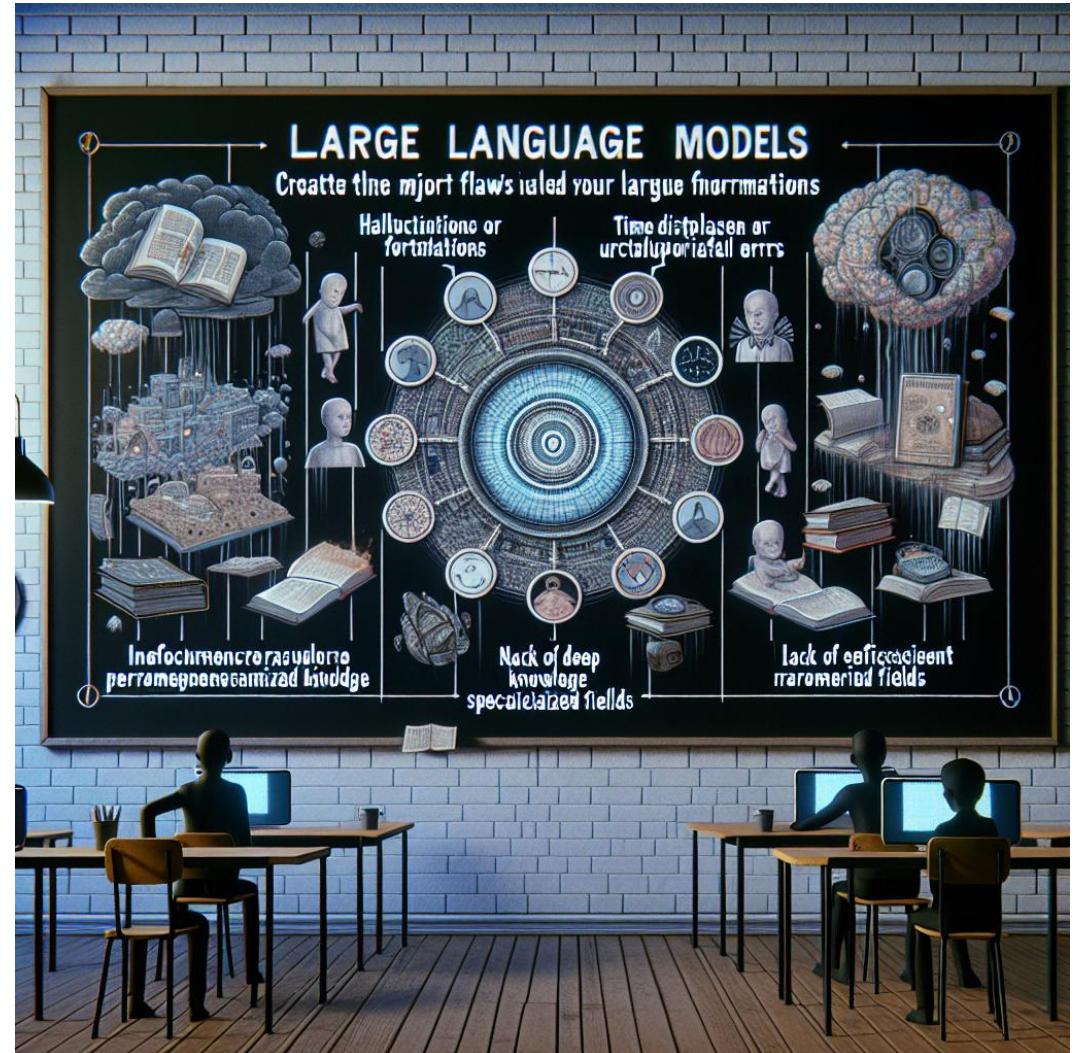
► Background

Drawbacks of LLMs

- Hallucination
- Outdated information
- Low efficiency in parameterizing knowledge
- Lack of in-depth knowledge in specialized domains
- Weak inferential capabilities

Practical Requirements of Application

- Domain-specific accurate answering
- Frequent updates of data
- Traceability and explainability of generated content
- Controllable Cost
- Privacy protection of data



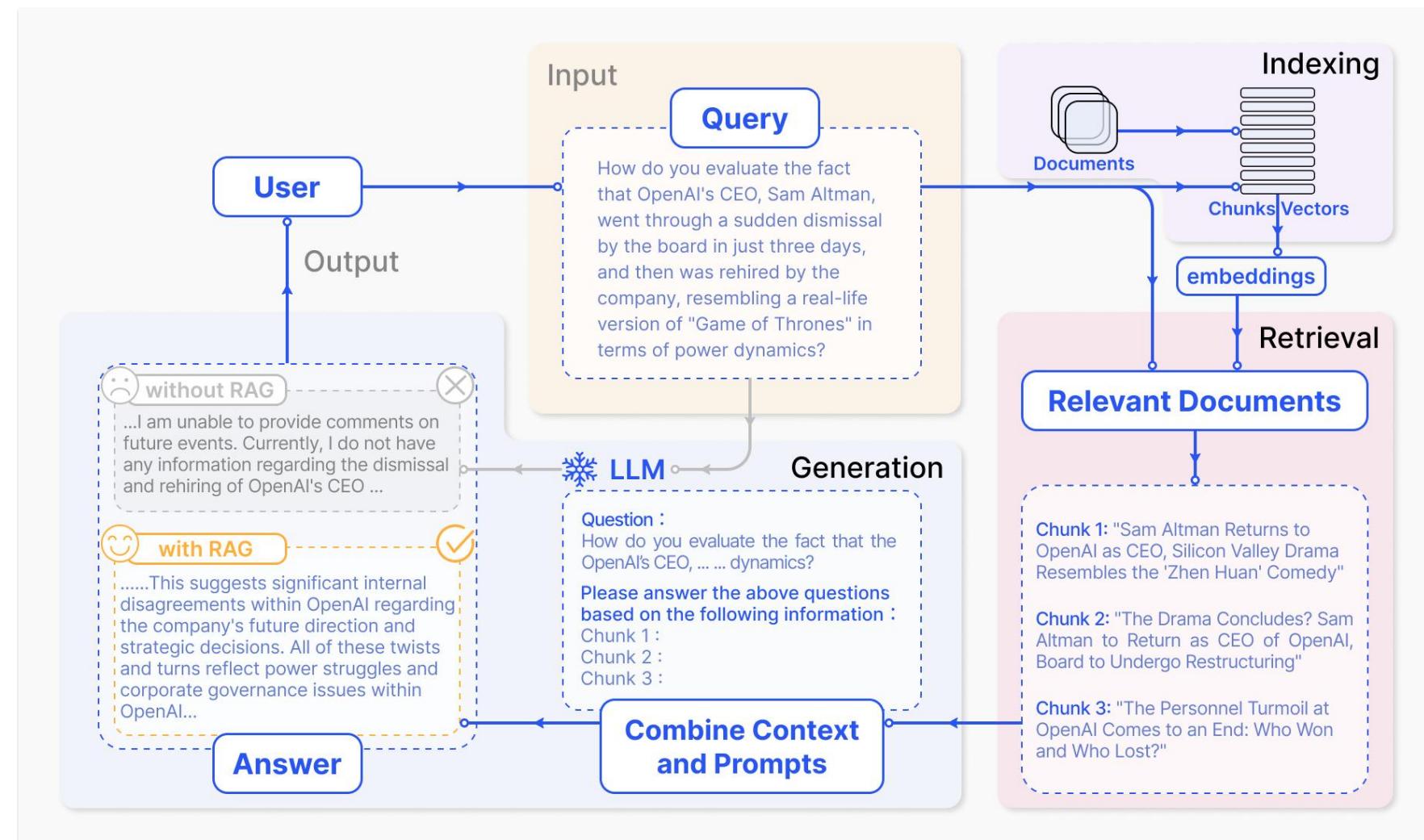
Draw by DALL-E-3

► Retrieval-Augmented Generation (RAG)

When answering questions or generating text, it first **retrieves relevant information** from a large number of documents, and then LLMs generates answers based on this information.

By attaching a **external knowledge base**, there is no need to retrain the entire large model for each specific task.

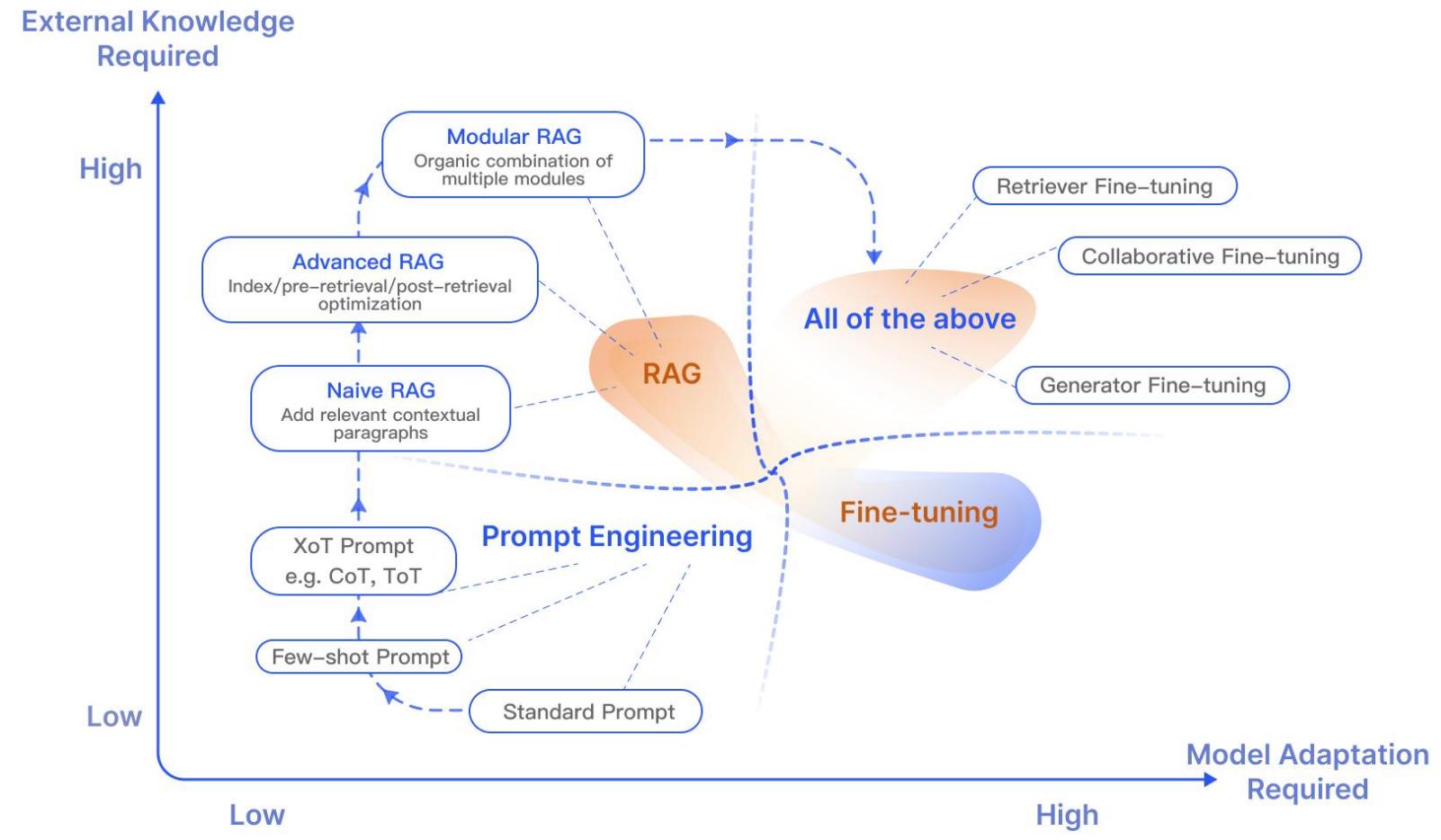
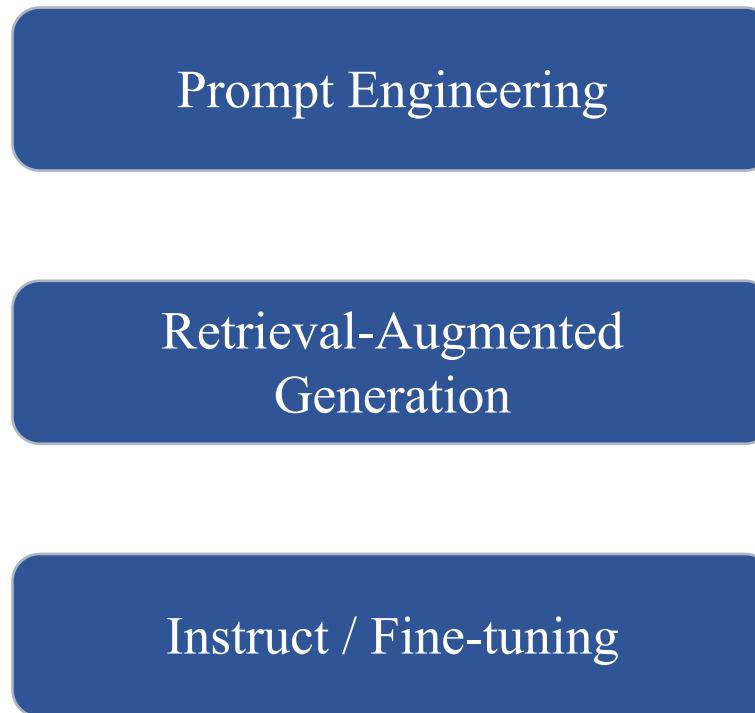
The RAG model is especially suitable for **knowledge-intensive** tasks.



A typical case of RAG

► Symbolic Knowledge or Parametric Knowledge

Ways to optimize LLMs.



A typical case of RAG

RAG vs Fine-tuning

Feature Comparison	RAG	Fine-Tuning
Knowledge Updates	Directly updating the retrieval knowledge base ensures that the information remains current without the need for frequent retraining, making it well-suited for dynamic data environments.	Stores static data, requiring retraining for knowledge and data updates.
External Knowledge	Proficient in leveraging external resources, particularly suitable for accessing documents or other structured/unstructured databases.	Can be utilized to align the externally acquired knowledge from pretraining with large language models, but may be less practical for frequently changing data sources.
Data Processing	Involves minimal data processing and handling.	Depends on the creation of high-quality datasets, and limited datasets may not result in significant performance improvements.
Model Customization	Focuses on information retrieval and integrating external knowledge but may not fully customize model behavior or writing style.	Allows adjustments of LLM behavior, writing style, or specific domain knowledge based on specific tones or terms.
Interpretability	Responses can be traced back to specific data sources, providing higher interpretability and traceability.	Similar to a black box, it is not always clear why the model reacts a certain way, resulting in relatively lower interpretability.
Computational Resources	Depends on computational resources to support retrieval strategies and technologies related to databases. Additionally, it requires the maintenance of external data source integration and updates.	The preparation and curation of high-quality training datasets, defining fine-tuning objectives, and providing corresponding computational resources are necessary.
Latency Requirements	Involves data retrieval, which may lead to higher latency.	LLM after fine-tuning can respond without retrieval, resulting in lower latency.
Reducing Hallucinations	Inherently less prone to hallucinations as each answer is grounded in retrieved evidence.	Can help reduce hallucinations by training the model based on specific domain data but may still exhibit hallucinations when faced with unfamiliar input.
Ethical and Privacy Issues	Ethical and privacy concerns arise from the storage and retrieval of text from external databases.	Ethical and privacy concerns may arise due to sensitive content in the training data.

► RAG Applications

Scenarios where RAG is applicable:

- Long-tail distribution of data
- Frequent knowledge updates
- Answers requiring verification and traceability
- Specialized domain knowledge
- Data privacy preservation

Q&A

RETRO (Borgeaud et al 2021)
REALM (Gu et al, 2020)
ATLAS (Izacard et al, 2023)

Fact Checking

RAG (Lewis et al, 2020)
ATLAS (Izacard et al, 2022)
Evi. Generator (Asai et al, 2022a)

Dialog

BlenderBot3 (Shuster et al.2022)
Internet-augmented generation (Komeili et a., 2022)

Summary

FLARE (Jiang et al, 2023)

Machine Translation

kNN-MT (Khandelwal et al., 2020)TRIME-MT (Zhong et al., 2022)

Code Generation

DocPrompting (Zhou et al., 2023)
Natural ProverWelleck et al., 2022)

Natural Language Inference

kNN-Prompt (Shi et al., 2022)
NPM (Min et al., 2023)

Sentiment analysis

kNN-Prompt (Shi et al., 2022)NPM (Min et al., 2023)

Commonsense reasoning

Raco (Yu et al, 2022)

PART 02

RAG Paradigms Shifting

► Naive RAG

Step1 Indexing

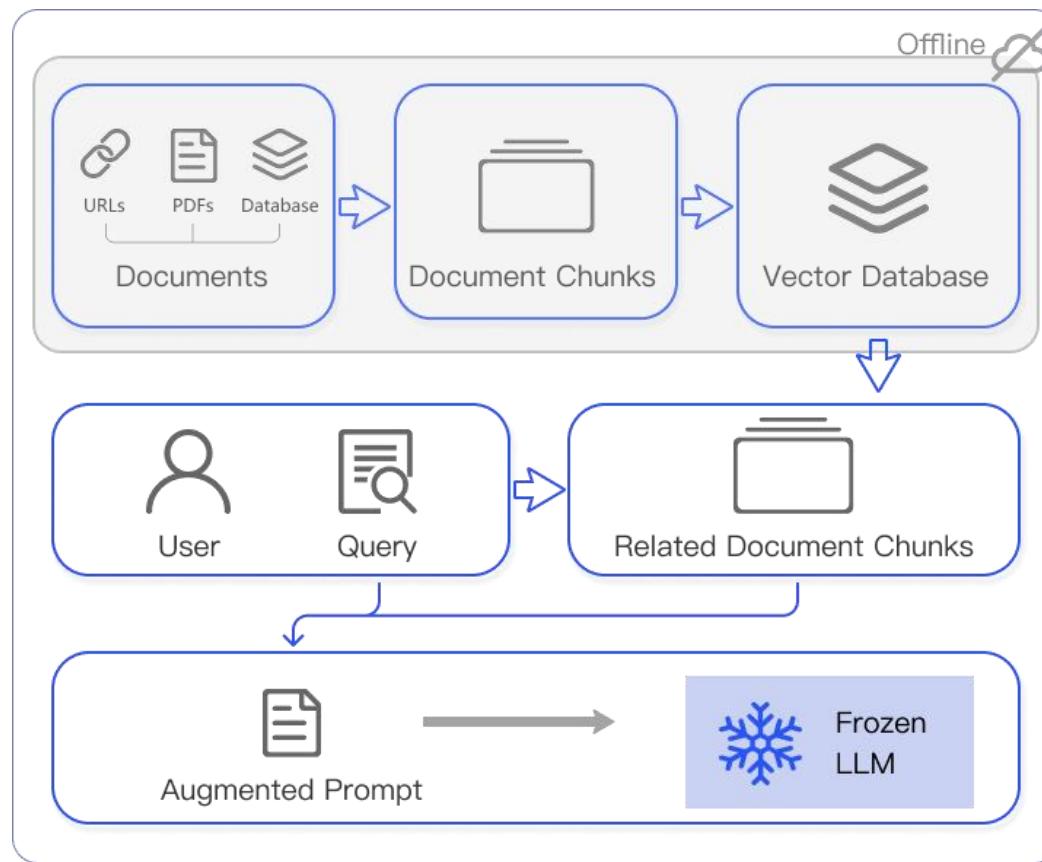
1. Divide the document into even chunks, each chunk being a piece of the original text.
2. Using the encoding model to generate an embedding for each chunk.
3. Store the Embedding of each block in the vector database.

Step2 Retrieval

Retrieve the k most relevant documents using vector similarity search.

Step3 Generation

The original query and the retrieved text are combined and input into a LLM to get the final answer



Naive RAG

Advanced RAG

Modular RAG

► Advanced RAG

Index Optimization → Pre-Retrieval Process → Retrieval →
Post-Retrieval Process → Generation

- **Optimizing Data Indexing:**

- sliding window, fine-grained

- segmentation、adding metadata

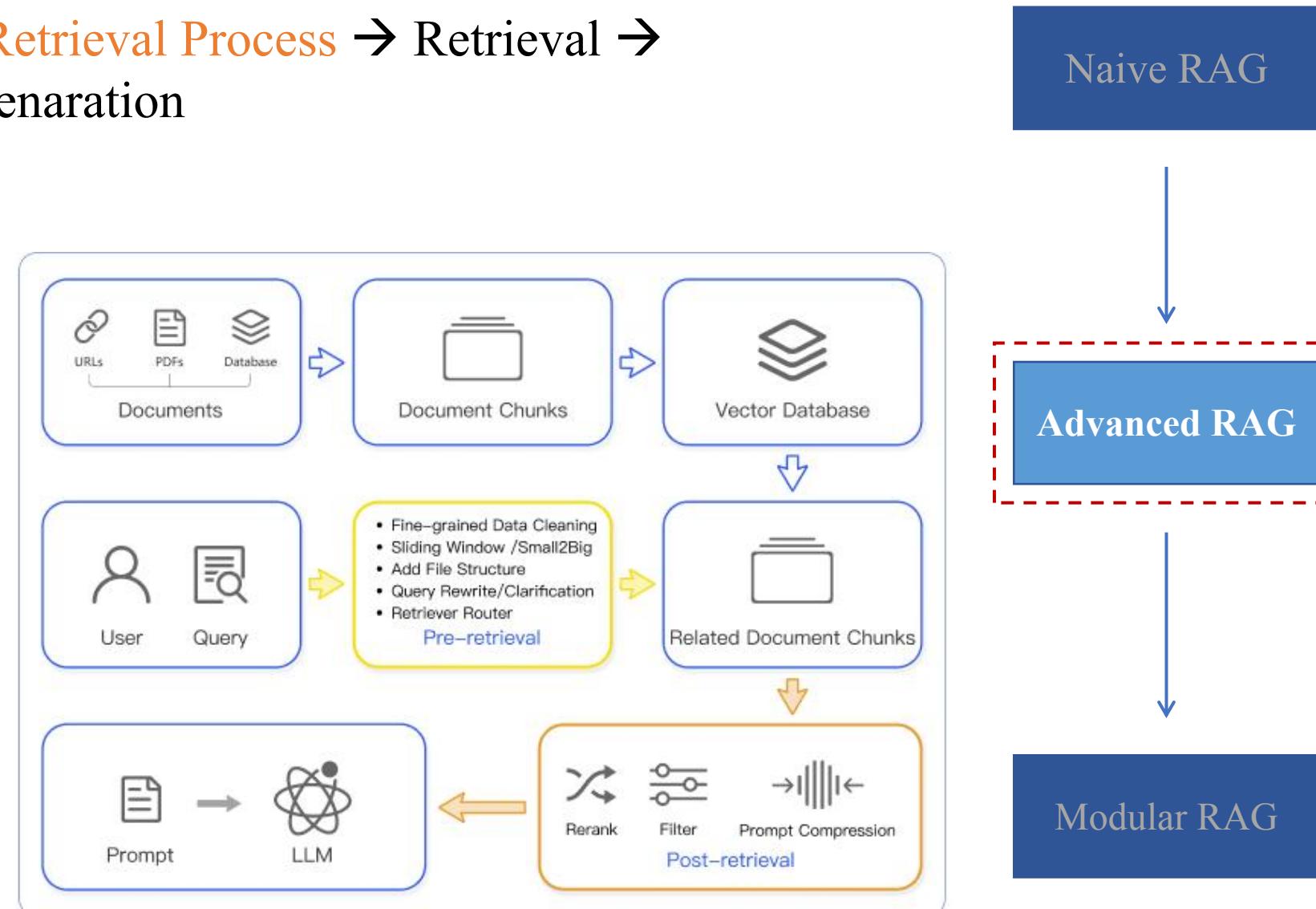
- **Pre-Retrieval Process:** retrieve

- routes, summaries, rewriting, and

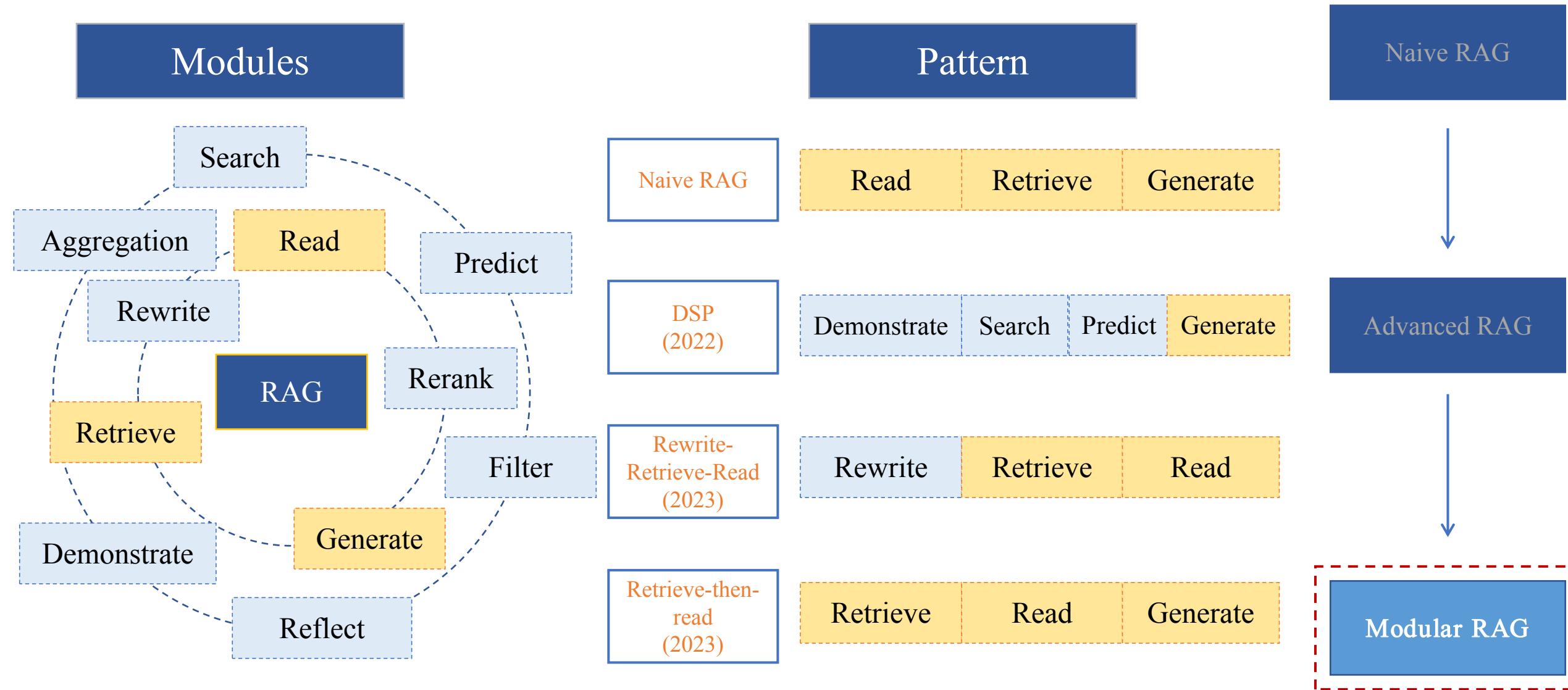
- confidence judgment

- **Post-Retrieval Process:** reorder,

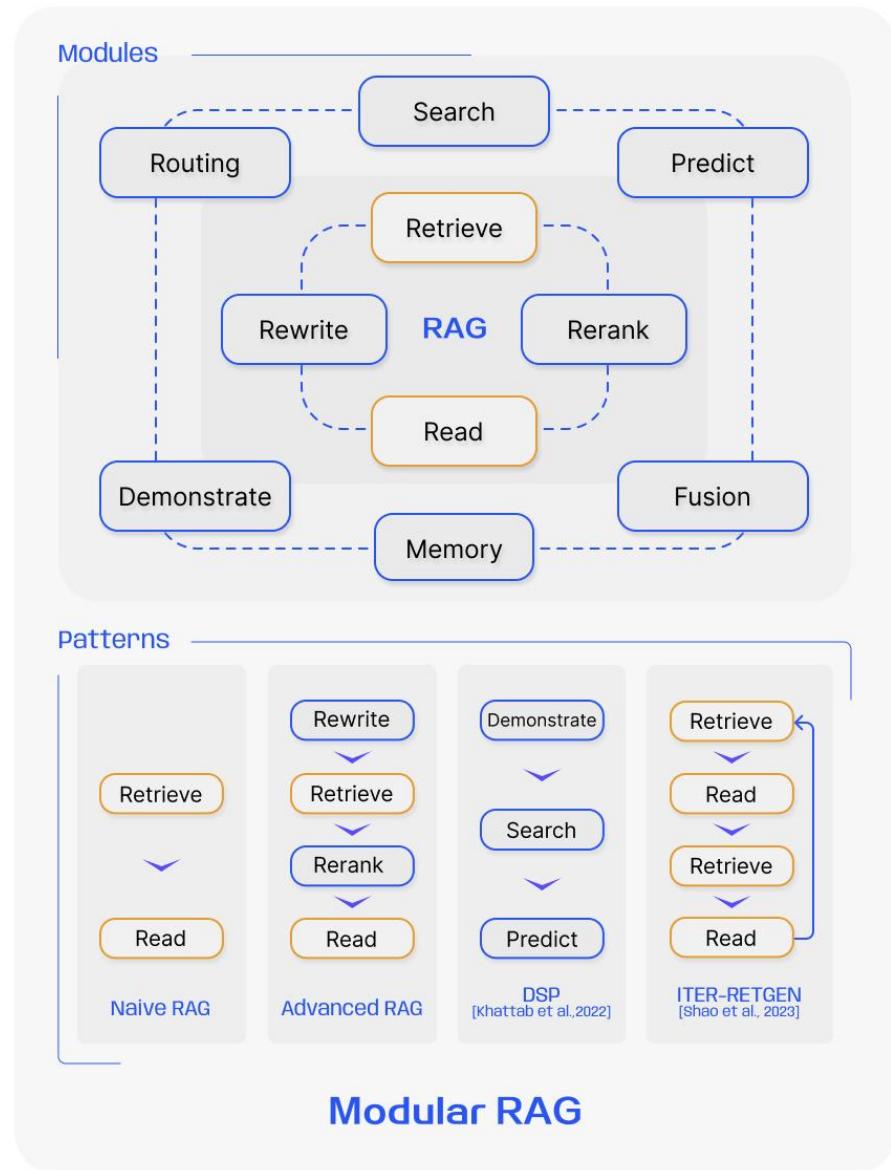
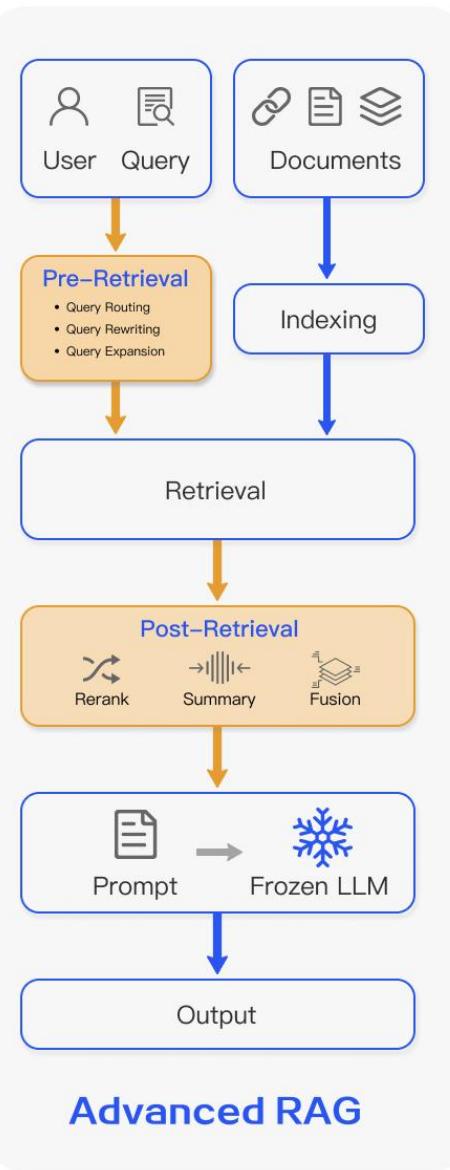
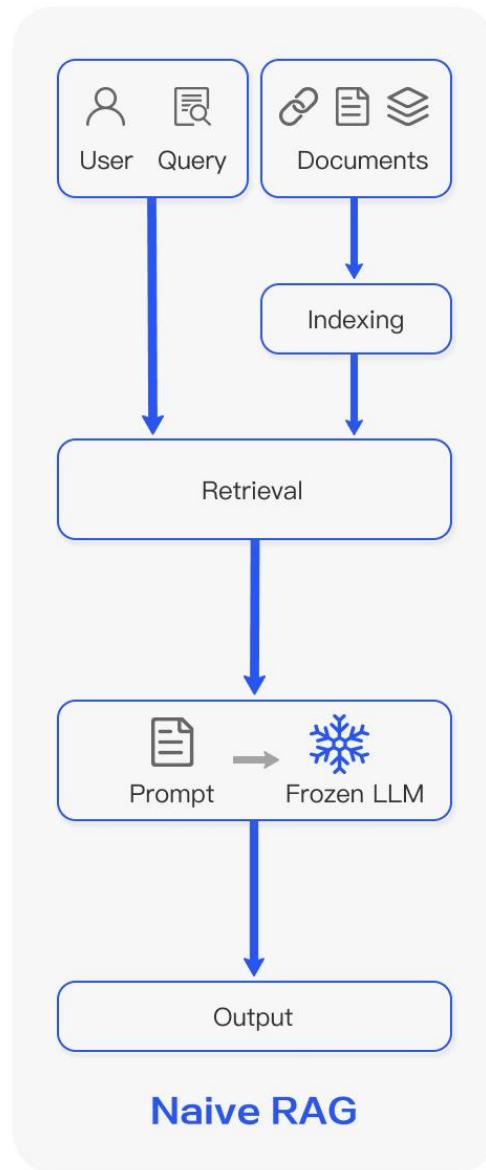
- filter content retrieval



► Modular RAG



► Comparison of RAG Paradigms



► The three key questions of RAG

What to retrieve ?

- Token
- Phrase
- Chunk
- Paragraph
- Entity
- Knowledge graph

When to retrieve ?

- Single search
- Each token
- Every N tokens
- Adaptive search

How to use the retrieved information ?

- Input/Data Layer
- Model/Intermediate Layer
- Output/Prediction Layer

Other Issues

Augmentation stage:

- Pre-training
- Fine-tuning
- Inference

Retrieval choice:

- BERT
- Roberta
- BGE
-

Model Collaboration
↔
Scale selectionz

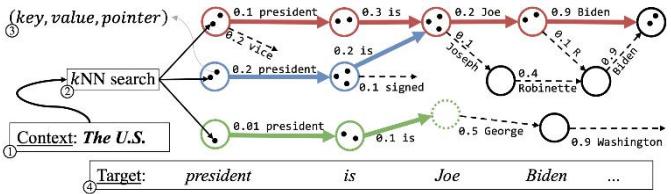
Generation choice:

- GPT
- Llama
- T5
-

► Key issue of RAG — What to retrieve

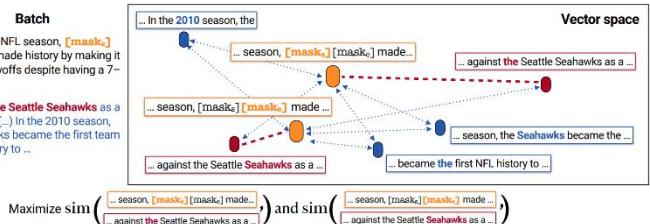
coarse
↑
Retrieval granularity
meticulous
↓
low

Chunk | In-Context RAG 2023

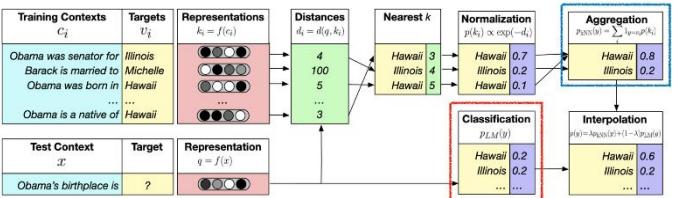


The search is **broad**, recalling a large amount of information, but with low **accuracy**, high coverage but includes much **redundant information**.

Phrase | NPM 2023



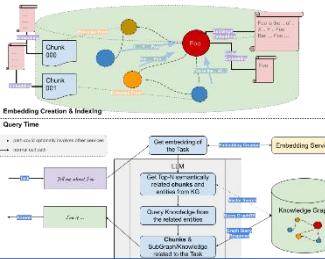
Token | KNN-LMM 2019



It excels in handling **long-tail** and cross-domain issues with **high computational efficiency**, but it requires **significant storage**.

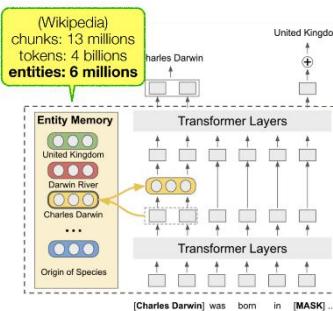
level of structuration

Knowledge Graph | 2023



Richer semantic and **structured information**, but the retrieval efficiency is lower and is limited by the quality of KG.

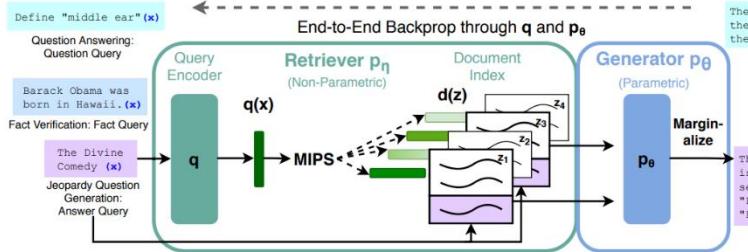
Entity | EasE 2022



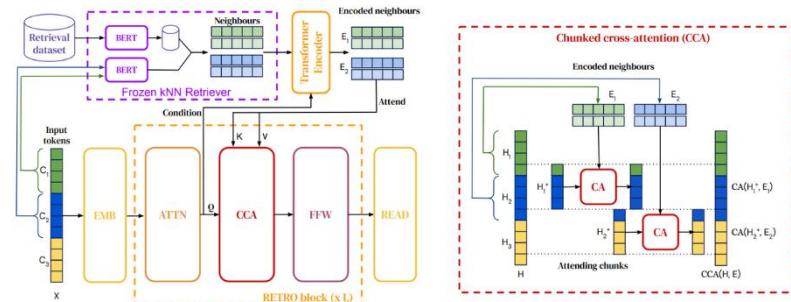
High

► Key issue of RAG — How to use the retrieved content

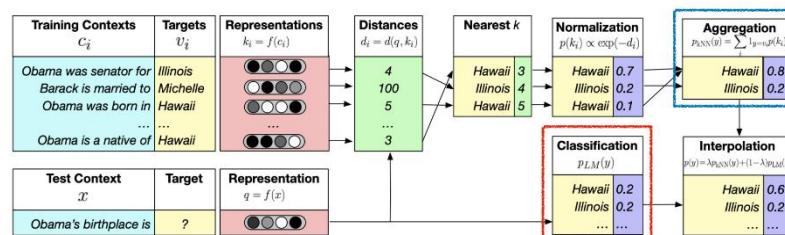
Integrating the retrieved information into different layers of the generation model, during inference process.



Input / Date layer



Model / Interlayer



Output /Prediction layer

Using simple, but unable to support the retrieval of more knowledge blocks, and the optimization space is limited.

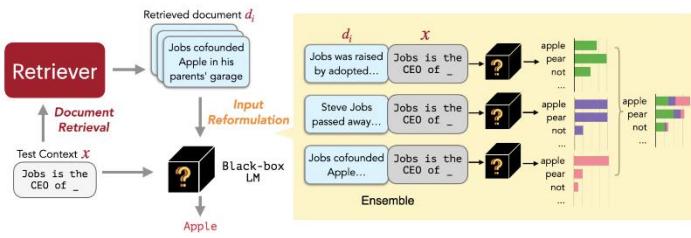
Supports the retrieval of more knowledge blocks, but introduces additional complexity and must be trained.

Ensuring the output results are highly relevant to the retrieval content, but the efficiency is low.

► Key issue of RAG — When to retrieve

High efficiency, but low relevance of the retrieved documents

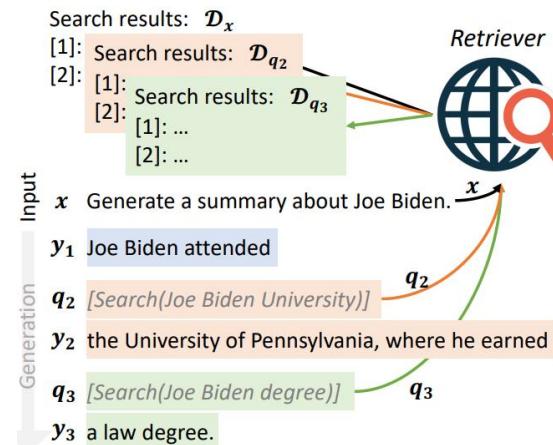
Once | Replug 2023



Conducting once search during the reasoning process.

Balancing efficiency and information might not yield the optimal solution

Adaptive | Flare 2023



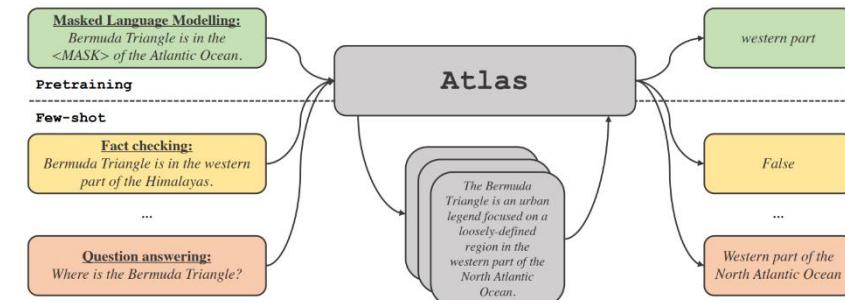
Adaptively conduct the search.

Low

Retrieval frequency

A large amount of information with low efficiency and redundant information.

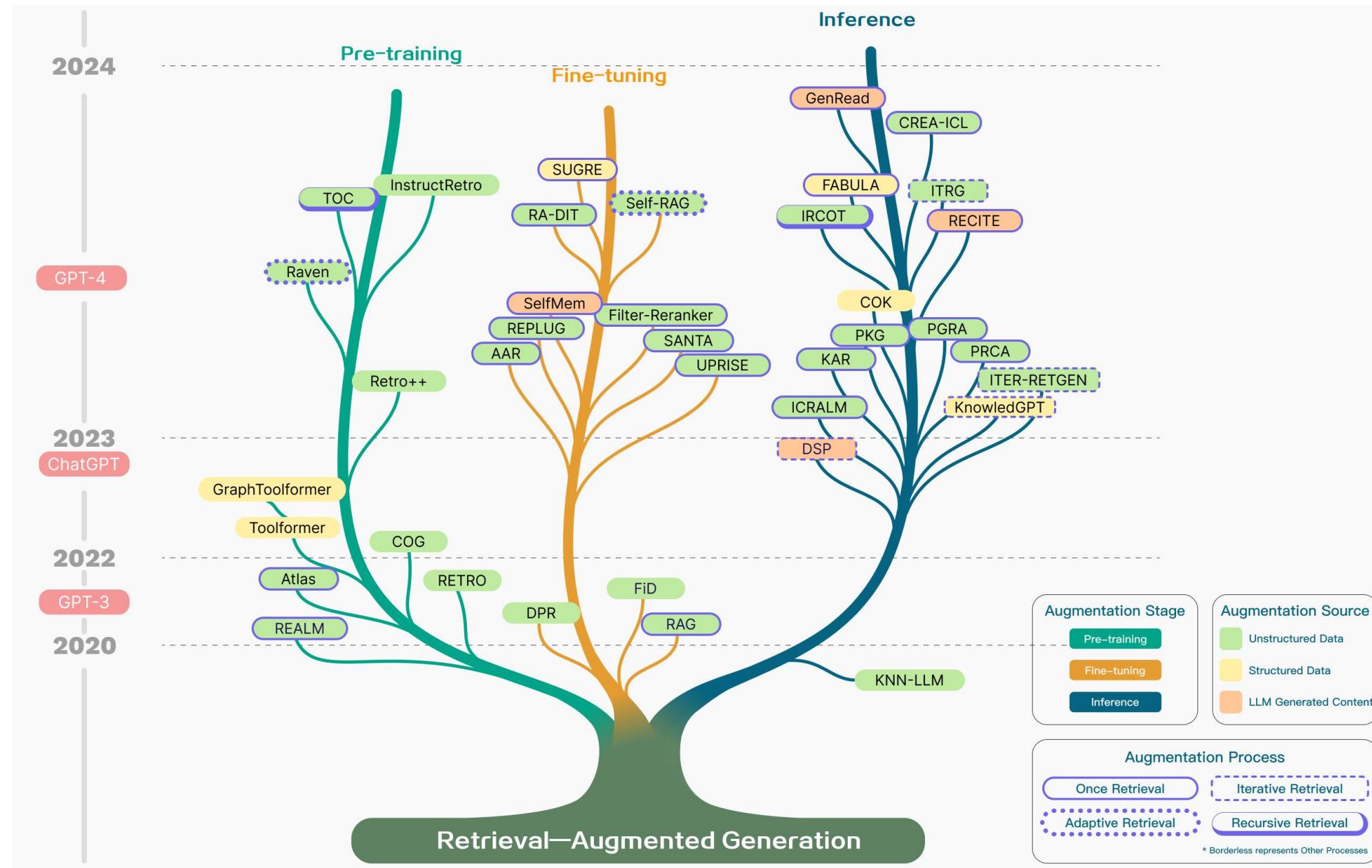
Every N Tokens | Atlas 2023



Retrieve once for every N tokens generated.

High

► Overview of RAG Development



PART 03

Key Technologies and Evaluation

► Techniques for Better RAG —— Data indexing optimization

Chunk Optimization

Small-2-Big

Embedding at sentence level expand the window during generation process.

Slidingwindow

Sliding chunk covers the entire text, avoiding semantic ambiguity

Summary

Retrieve documents through summaries, then retrieve text blocks from the documents.

Adding Metadata

Example

Page

Time

Type

Document Title

Metadata Filtering/Enrichment

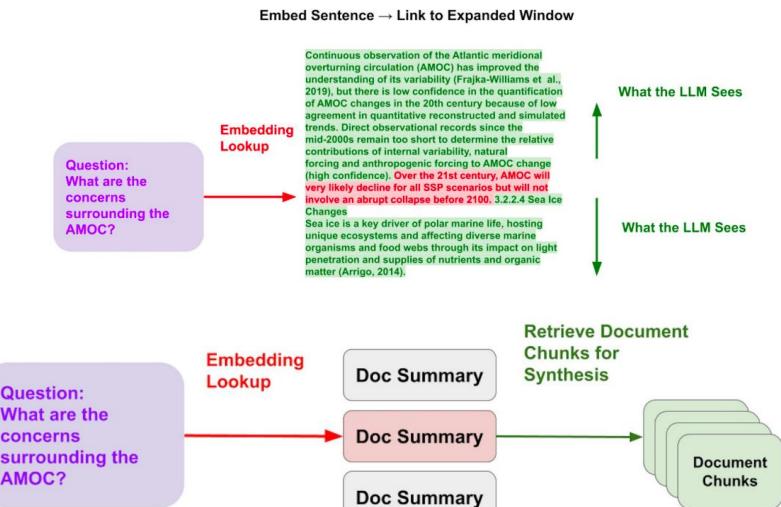
Pseudo Metadata Generation

Enhance retrieval by generating a hypothetical document for the incoming query and creating questions that the text block can answer.

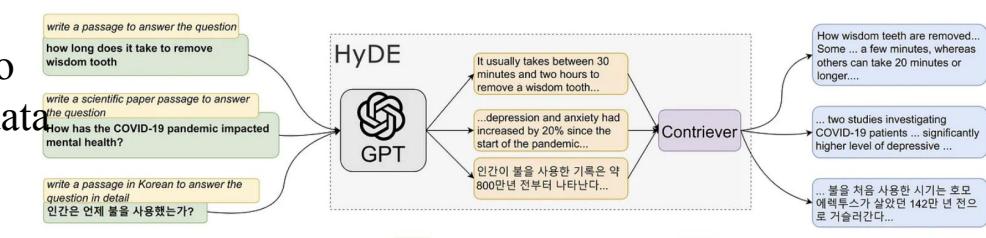
Metadata filter

Dissect and annotate the document. During the query, infer metadata filters in addition to semantic queries

Small-2-Big

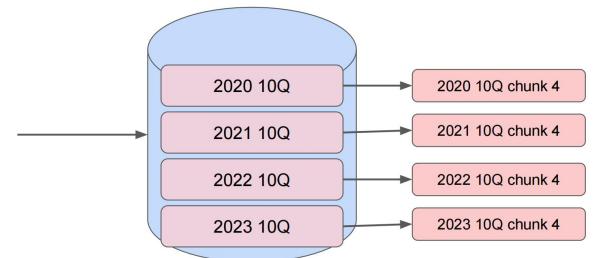


Pseudo Metadata



Metadata filter

query_str:
<query_embedding>
Metadata tags:
<metadata_tags>

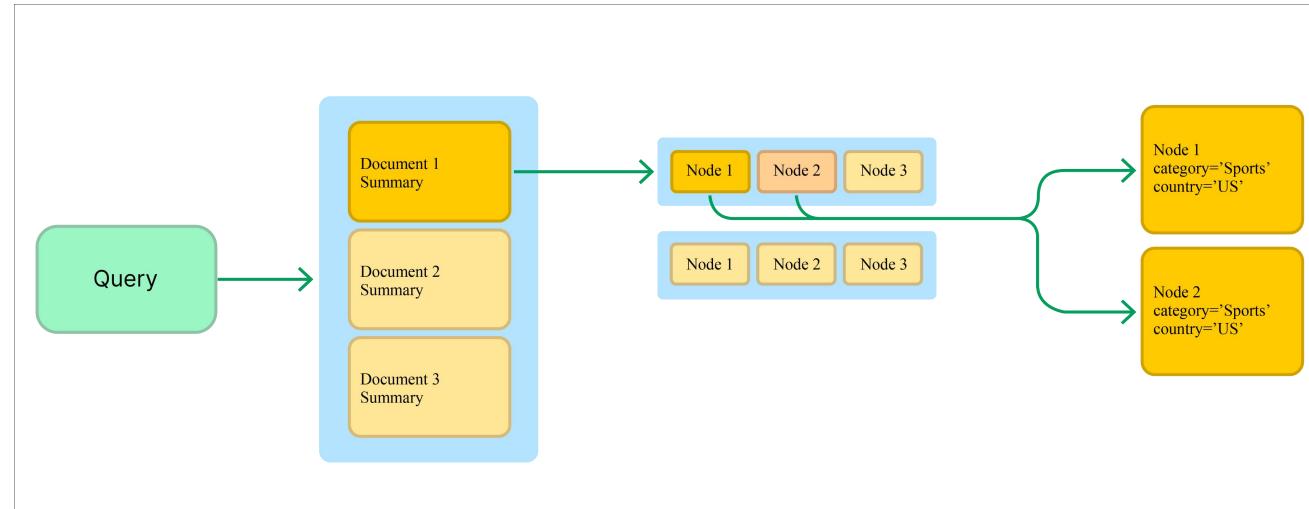


► Techniques for Better RAG — Structured Corpus

Hierarchical Organization of Retrieval Corpora

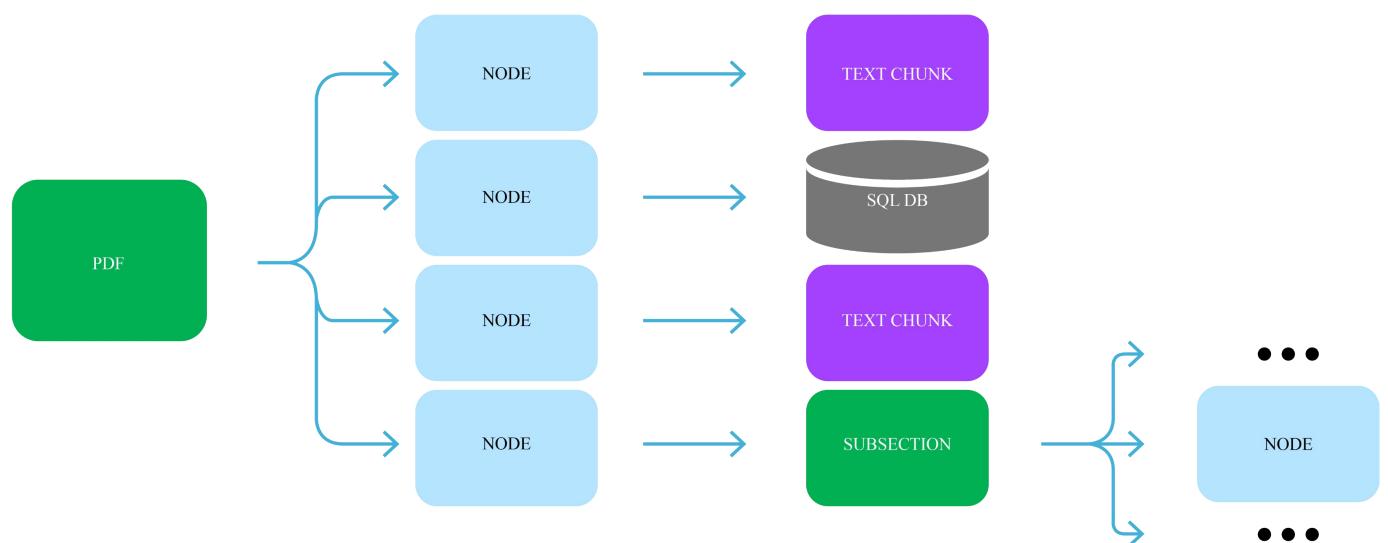
- Summary → Document

Replace document retrieval with summary retrieval, not only retrieving the most directly relevant nodes, but also exploring additional nodes associated with those nodes.



- Document → Embedded Objects

Documents have embedded objects (such as tables, charts), first retrieve entity reference objects, then query underlying objects, such as document blocks, databases, sub-nodes.



► Techniques for Better RAG —— Retrieval Source Optimization



Phrases

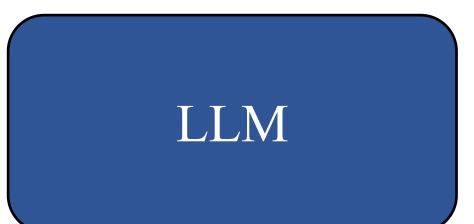
Prompt

Cross-linguistic



Triples

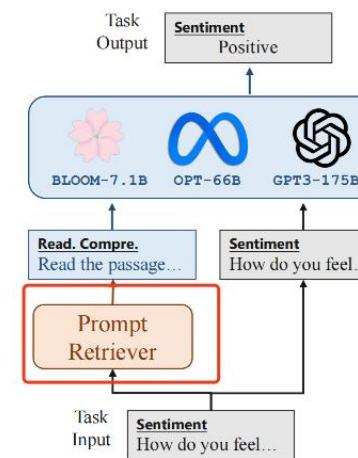
Subgraphs



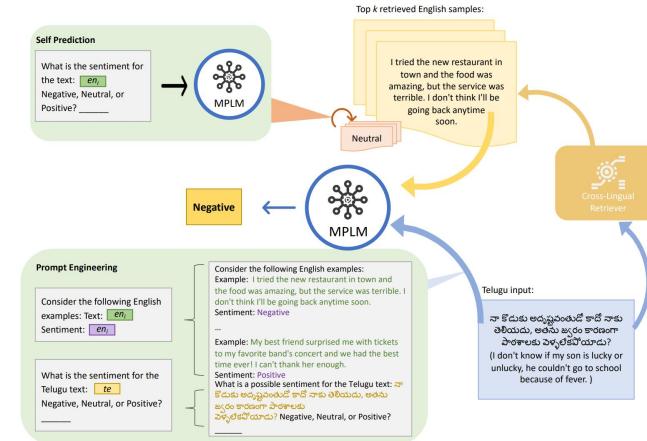
LLM Memory

Generated Text

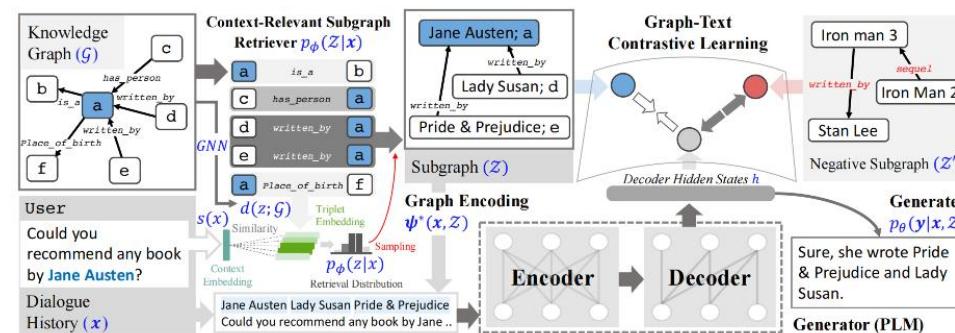
Generated Code



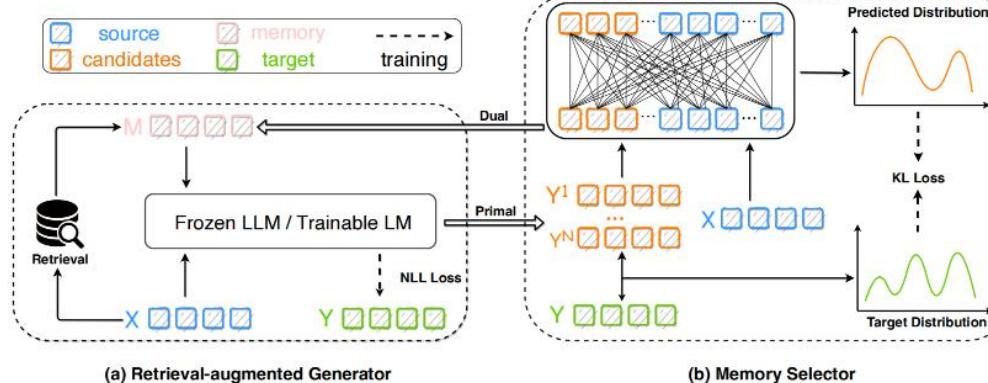
Prompt | UPRISE [Cheng et al., 2023]



Cross-language| CREA-ICL [Li et al., 2023]



Subgraph | SUGRE [Kang et al., 2023]



Memory | Selfmem [Cheng et al., 2023]

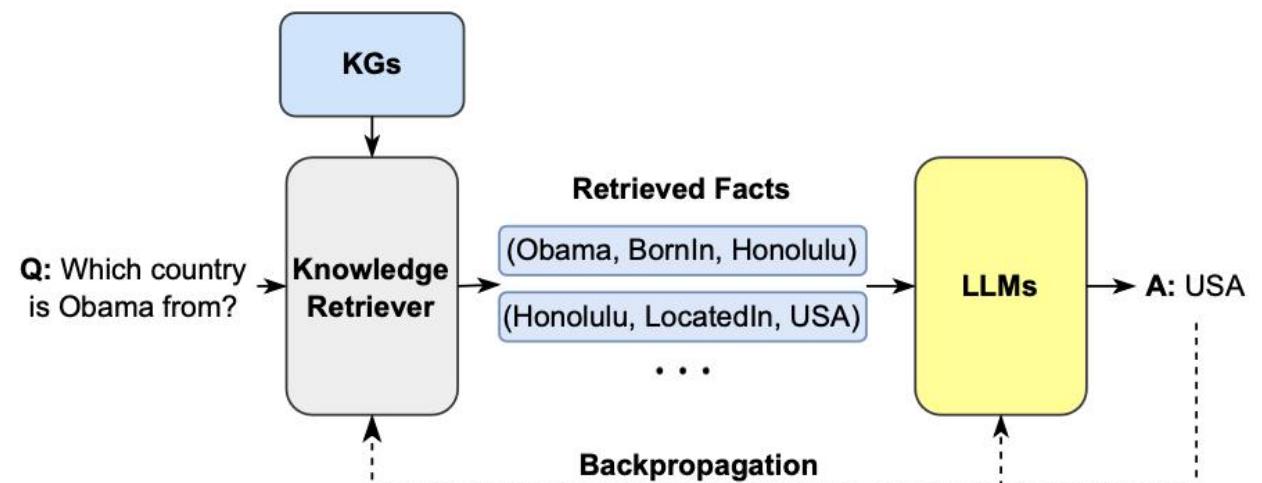
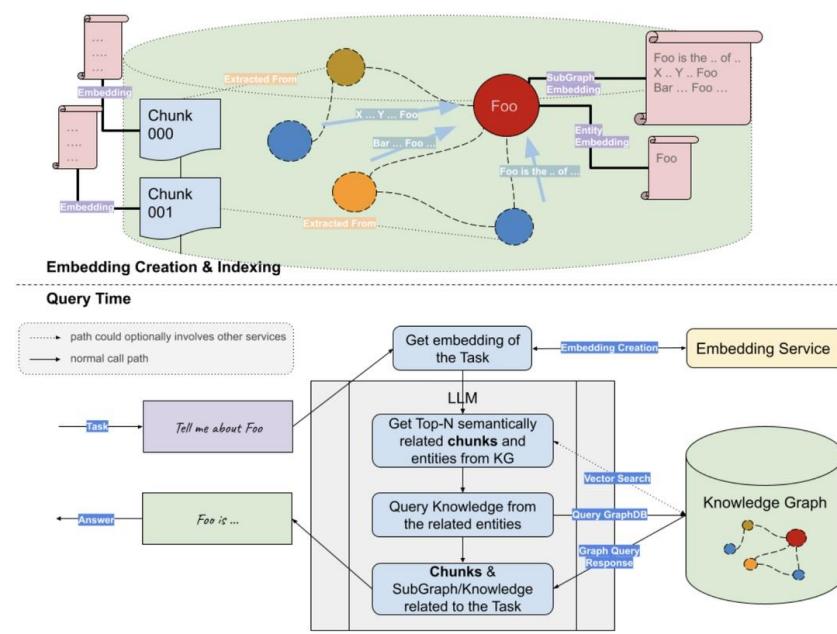
► Techniques for Better RAG — KG as a Retrieval Data Source

➤ GraphRAG

- Extract entities from the user's input query, then construct a subgraph to form context, and finally feed it into the large model for generation.

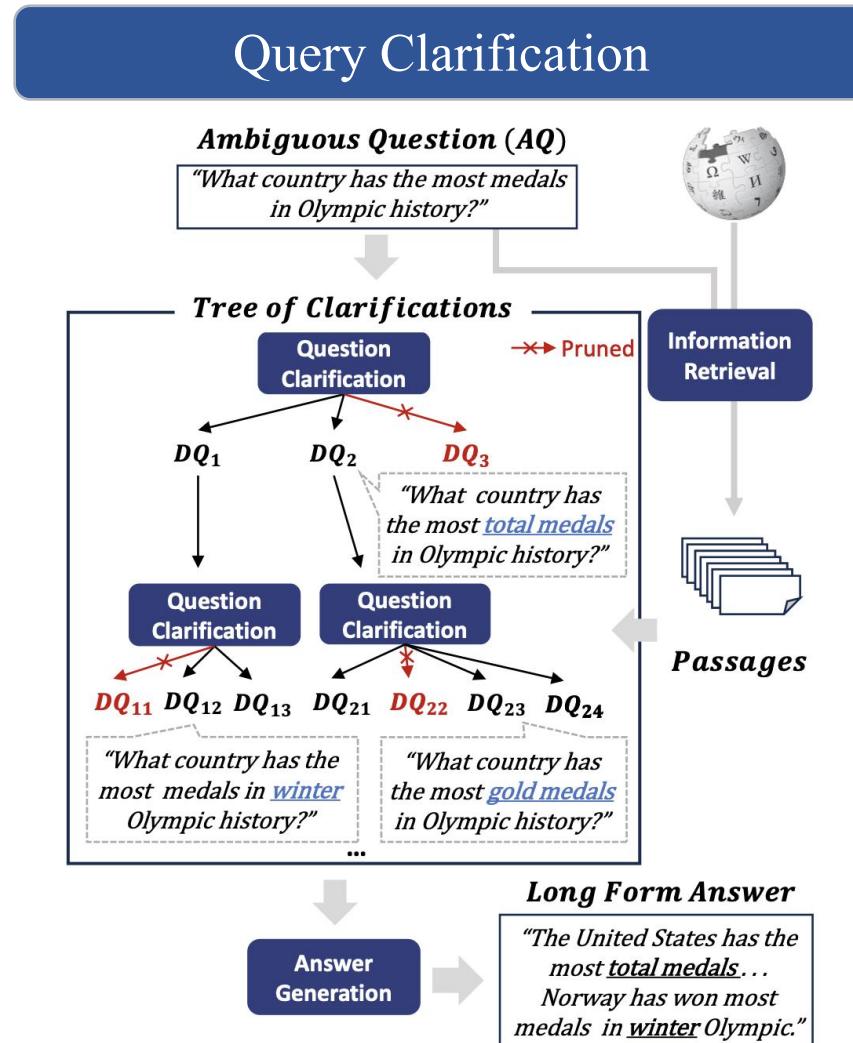
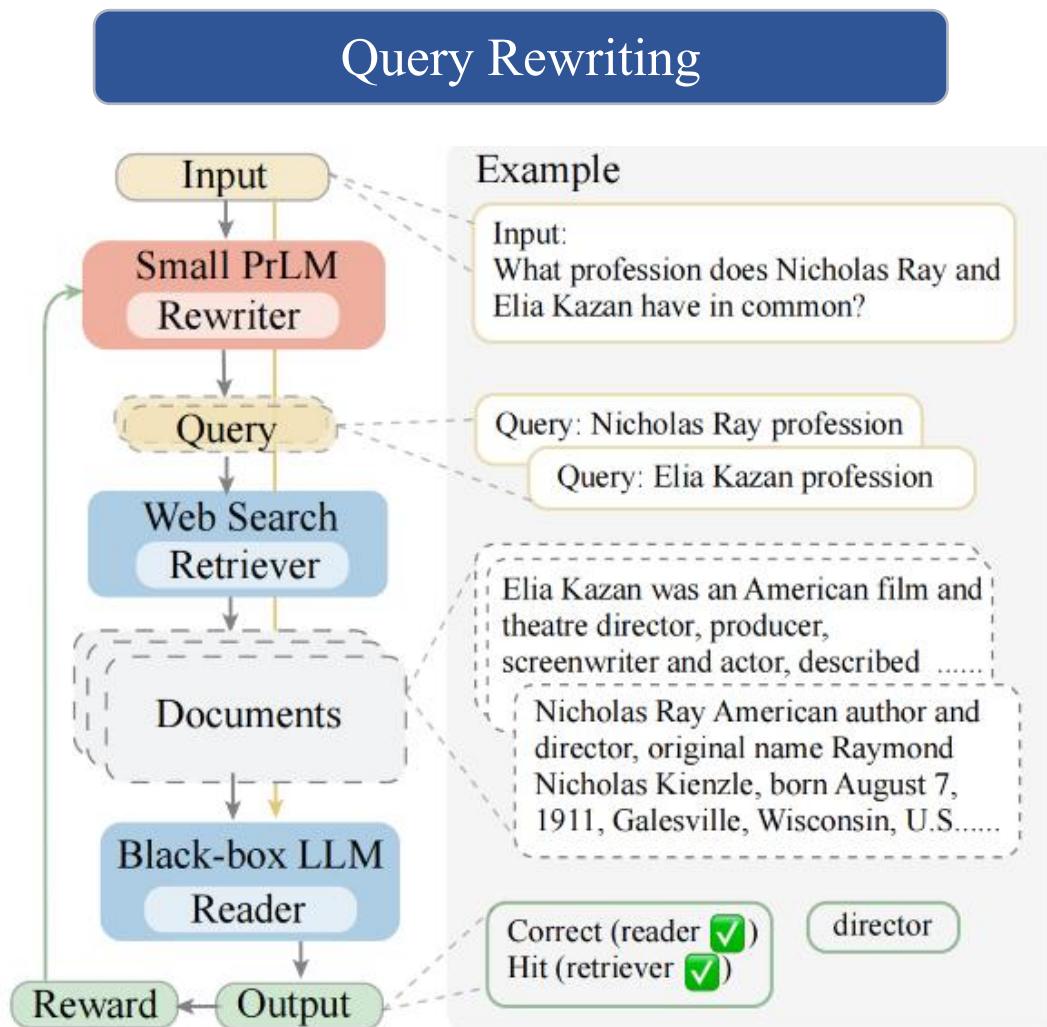
➤ Implementation

- Use LLM (or other models) to extract key entities from the question.
- Retrieve subgraphs based on entities, delving to a certain depth, such as 2 hops or even more.
- Utilize the obtained context to generate answers through LLM.



► Techniques for Better RAG —— Query Optimization

Questions and answers do not always possess high semantic similarity; adjusting the Query can yield better retrieval results.

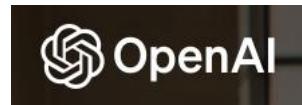


► Techniques for Better RAG —— Embedding Optimization

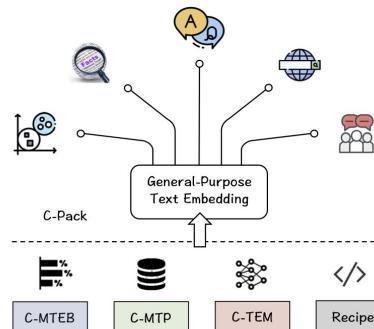
Selecting a More Suitable Embedding Provider



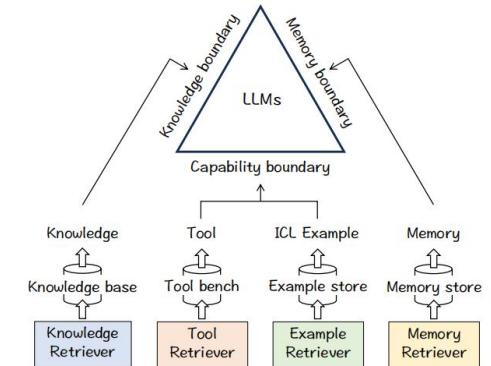
VOYAGE AI



BAAI
智源研究院

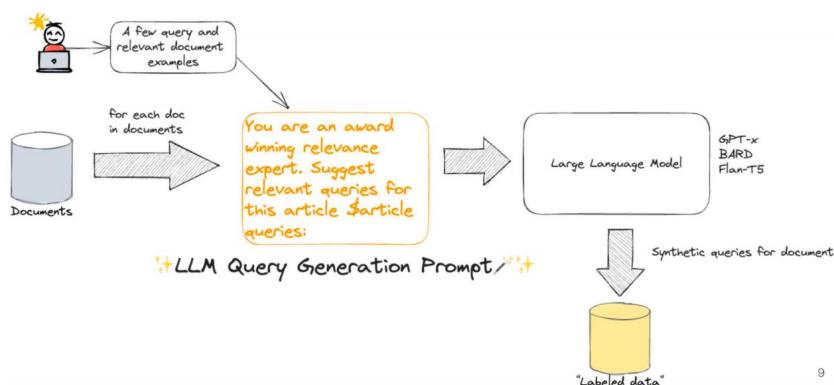


BAAI-General-Embedding (BGE)

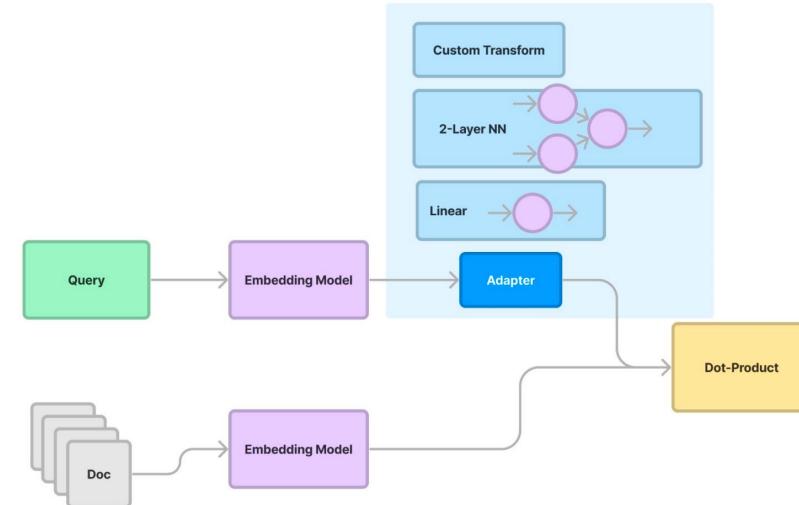


LLM-Embedder(BGE2) [Aksitov et al., 2023]

Fine-tuning the Embedding Model



Fine-tuning According to Domain-Specific
Repositories and Downstream Tasks



Fine-tuning the Adapter Module to Align the Embedding
Model with the Retrieval Repository

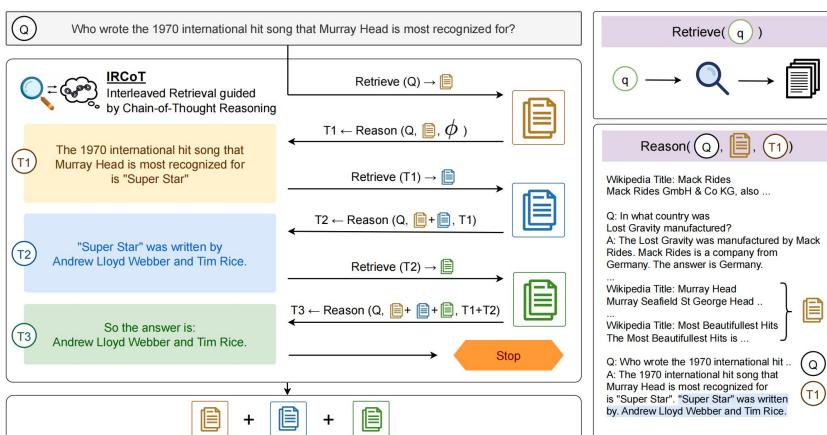
► Techniques for Better RAG — Retrieval Process Optimization

Iterative

Iteratively Retrieving from the Corpus to Acquire More Detailed and In-depth Knowledge



ITER [Feng et al., 2023]

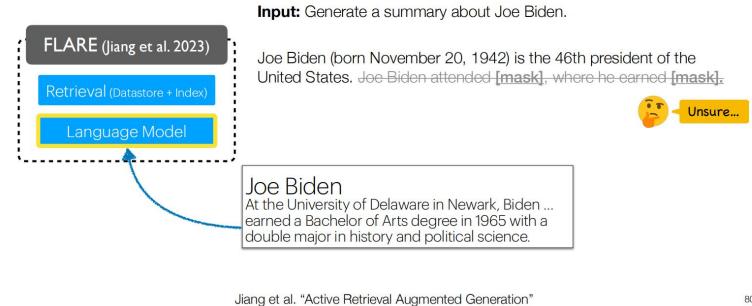


IRCOT [Trivedi et al., 2022]

Adaptive

Dynamically Determined by the LLM, the Timing and Scope of Retrieval

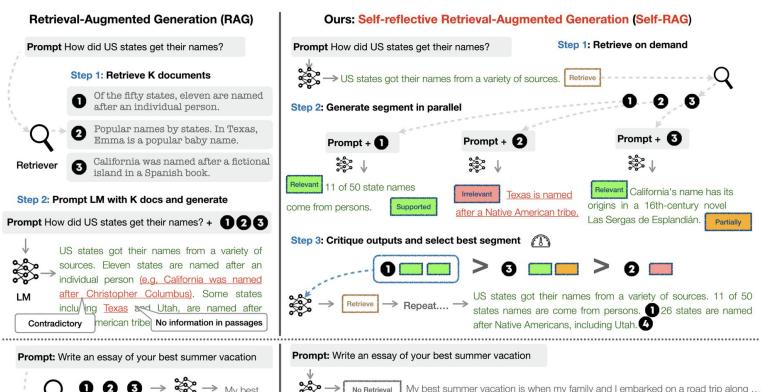
FLARE [Jiang et al., 2023]



Jiang et al. "Active Retrieval Augmented Generation"

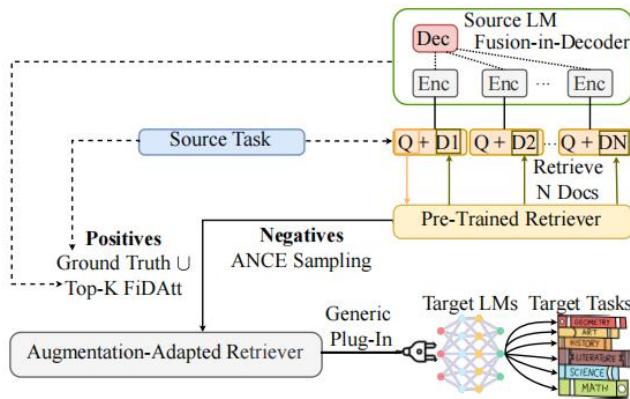
80

Self-RAG [Asai et al., 2023]



► Techniques for Better RAG — Hybrid (RAG + Fine-tuning)

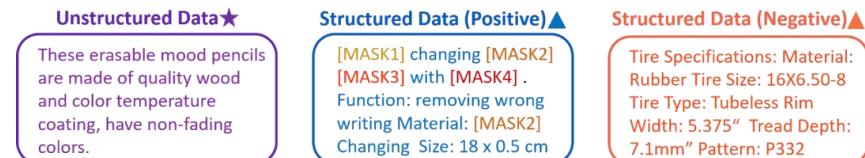
Retriever Fine-Tuning



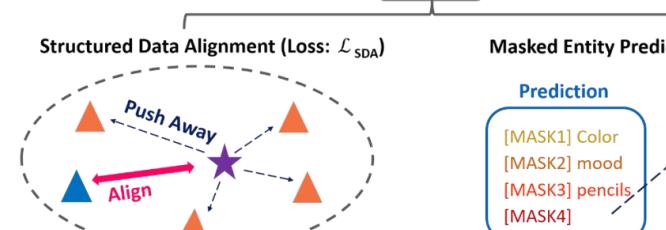
Highly Adaptive General-Purpose Retrieval Plugin

AAR [Yu et al., 2023]

Generator Fine-Tuning

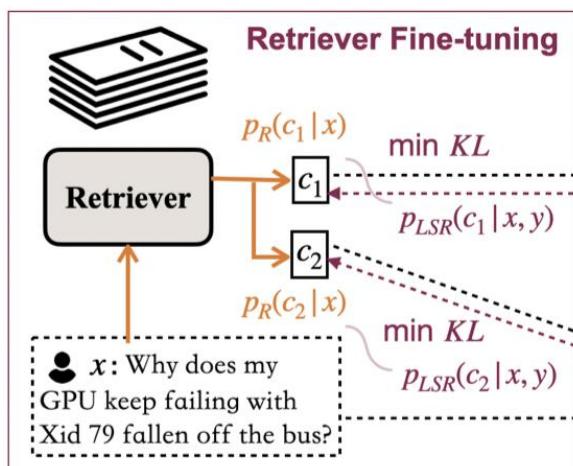


Augment with Structural Information Integration



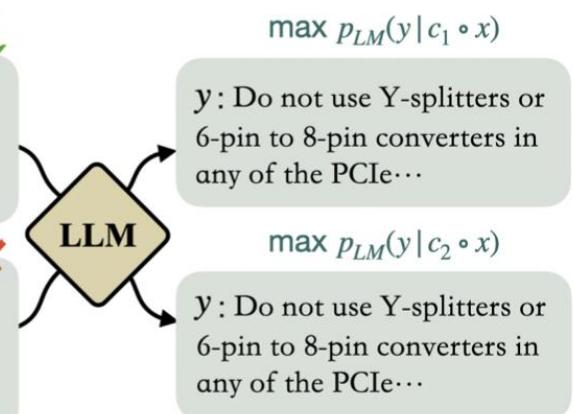
SANTA [Li et al., 2023]

Collaborative Fine-Tuning



Retrieval-augmented Instruction Tuning

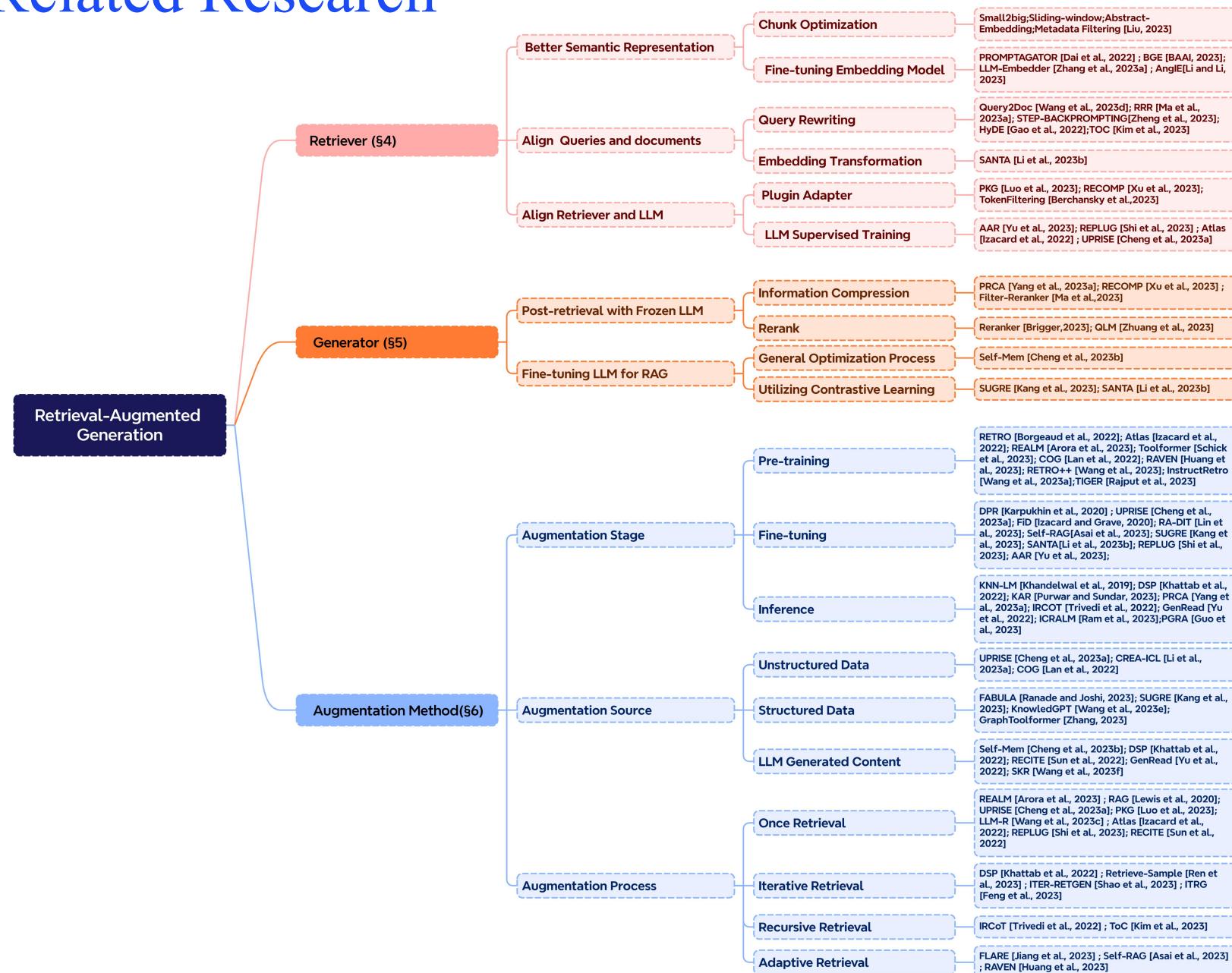
- 1 **Background:** I assume that the BGA chip has damage to the substrate level
... \n\n**Q:** Why does my GPU keep failing with Xid 79 fallen off the bus? **A:**
- 2 **Background:** Microsoft should withdraw from the hardware market ... \n\n**Question:** Why does my GPU keep failing with Xid 79 fallen off the bus? **Answer:**



RA-DIT [Lin et al., 2023]

- **R-FT**
Minimizing the KL Divergence Between the Retriever Distribution and LLM Preferences
- **LM-FT**
Maximizing the Likelihood of the Correct Answer Given Retrieval-Augmented Instructions

► Summary of Related Research



► How to Evaluate the Effectiveness of RAG

Evaluation Methods

Independent Evaluation

Retriever

Evaluate the Quality of Text Blocks Retrieved by the Query Metrics: MRP, Hit Rate, NDCG

Generation/Synthesis

Quality of Context Enhanced with Retrieved Documents Evaluation Metrics: Context Relevance

End-to-End Evaluation

Evaluate the content ultimately generated by the model.

By generated content

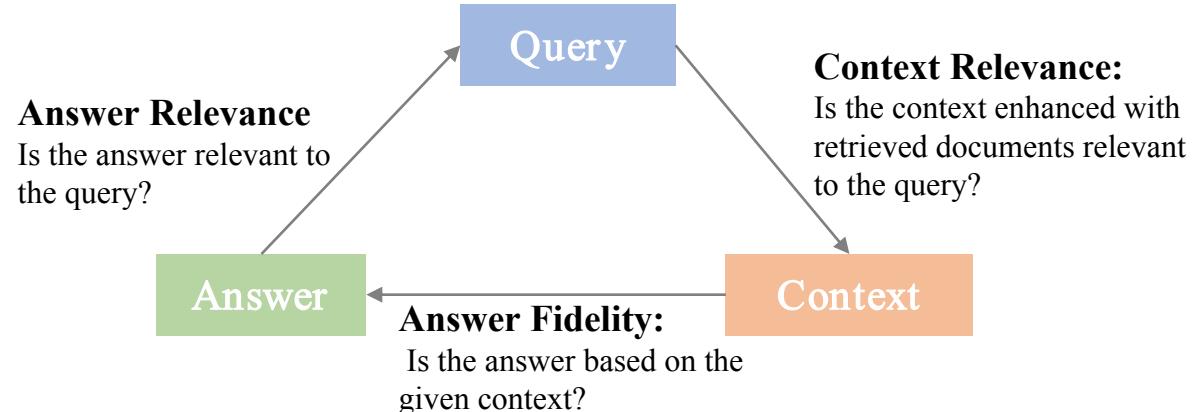
With labels: EM, Accuracy
Without labels: Fidelity, Relevance, Harmlessness

By evaluation method

Human evaluation
Automatic evaluation (LLM judge)

Key Metrics & Capabilities

Key Metrics



Assessment Framework

Use LLM as the adjudicator judge.

TruLens

RAGAS

ARES

Based on handwritten prompt

Synthetic dataset + Fine-tuning + Ranking using confidence intervals

Evaluation

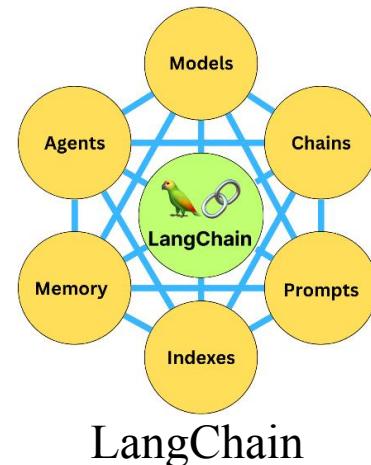
- Answer Fidelity
- Answer Relevance
- Contextual Relevance

PART 04

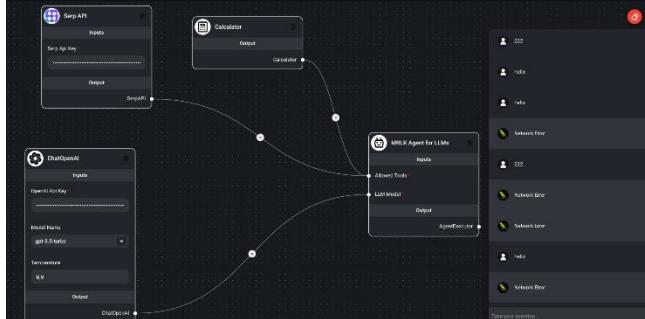
RAG Stack and Industry Practices

► Existing Tech Stack for RAG

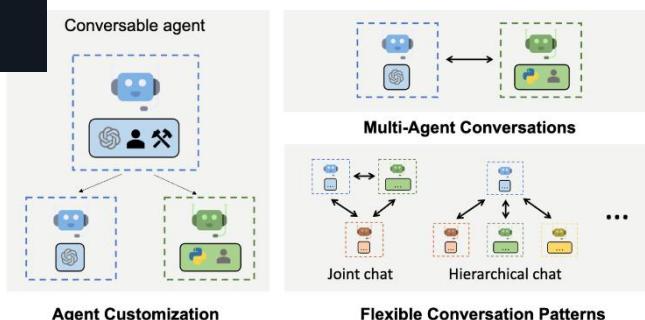
Name	Pros	Cons
LangChain	Modular, full-featured	Inconsistent behavior ,API conceals details, complexity and low flexibility.
LlamaIndex	Focus on RAG	Requires combination use, low customization.
FlowiseAI	Easy to get started, visualized workflows.	Does not support complex scenarios.
AutoGen	Adapts to multi-agent scenarios.	Low efficiency, requires multiple rounds of dialogue.



LlamaIndex



FlowiseAI

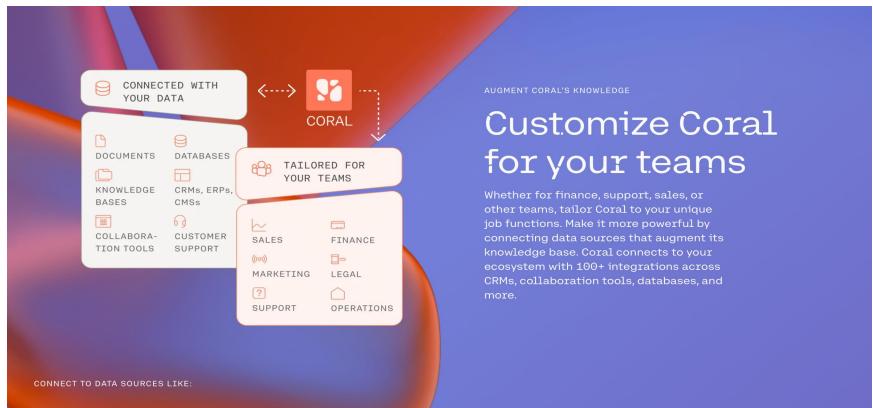


AutoGen

► RAG Industry Application Practices



NetEase - ChatBI



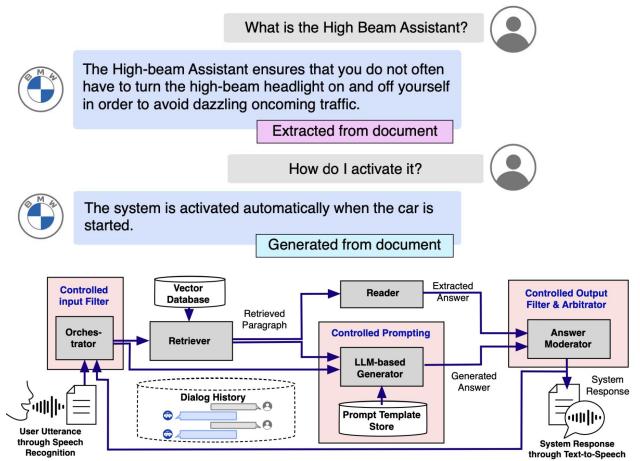
Cohere - Coral

The intelligent upgrade of traditional industries

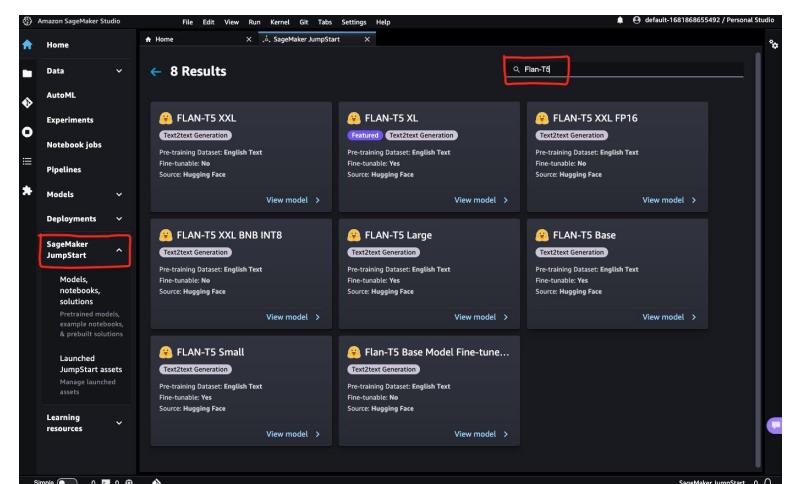


RAG

AI Toolchain
Enhancement



BMW - CarExpert

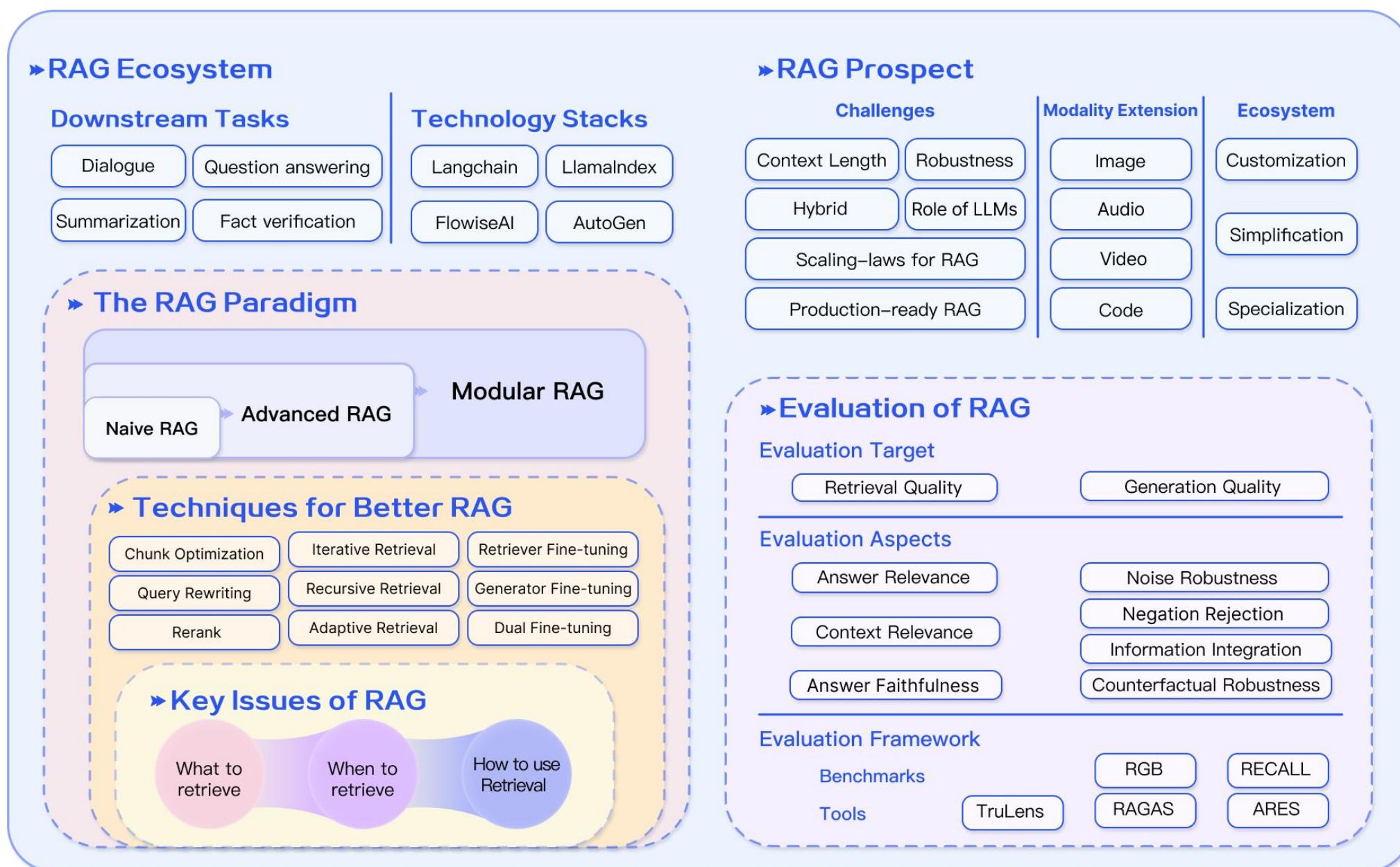


Amazon - Kendra

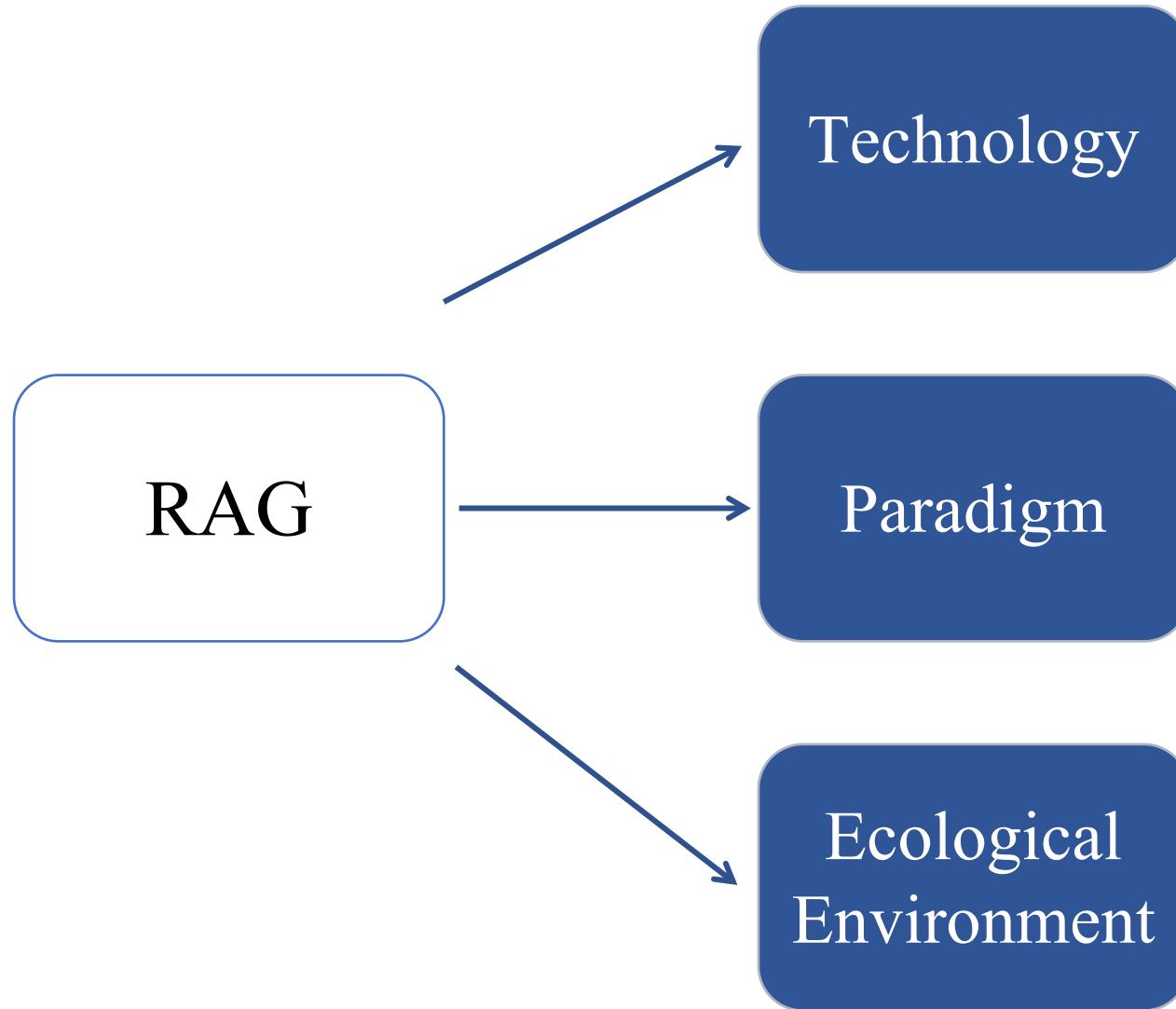
PART 06

Summary and Outlook

► Summary — The Framework of RAG



► Summary —— Three Trends of RAG



- The Scaling Law of RAG Models
- How to Improve the Efficiency of Retrieving Large-scale Data
- Mitigation of Forgetting in Long-context Scenarios
- Enhancement of Multimodal Retrieval
- Modularity Will Become Mainstream
- Patterns for Module Organization Await Refinement
- Evaluation Systems Need to Evolve and Improve with Time
- Preliminary Formation of Toolchain Technology Stack
- One-stop Platform Still Requires Polishing
- Explosion of Enterprise-level Applications

▶ Prospects — Existing Challenges of RAG

Further address the challenges faced by RAG itself

Long context

- Retrieved content is excessive, **exceeding window limit**.
- The context is too long to result **Lost in the Middle**.
- If the context **window is not limited**, is there still a need for RAG?

Coordination with FT

- How to simultaneously leverage the effects of **RAG** and **FT**.
- How do the two coordinate, how are they organized, is it in **Pipeline**, **alternating**, or **end-to-end**?

Robustness

- How to handle the **incorrect** content retrieved
- How to **filter** and **verify** the content retrieved.
- How to improve the model's **resistance to toxicity and noise**

Scaling Law

- Does the RAG model satisfy the **Scaling Law**
- Does RAG exhibit, or under what scenarios does it exhibit an **Inverse Scaling Law**

The role of LLMs

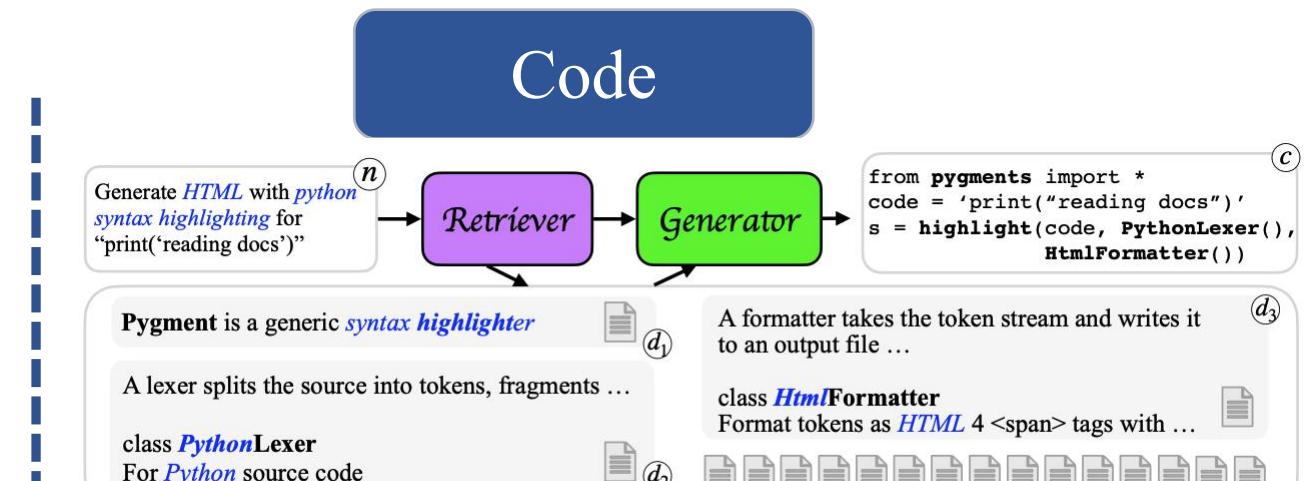
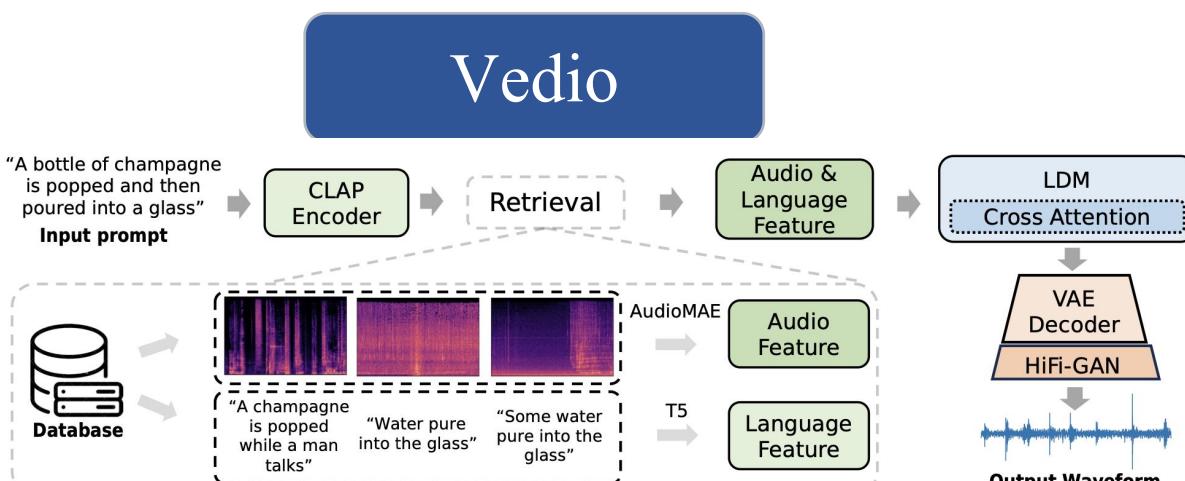
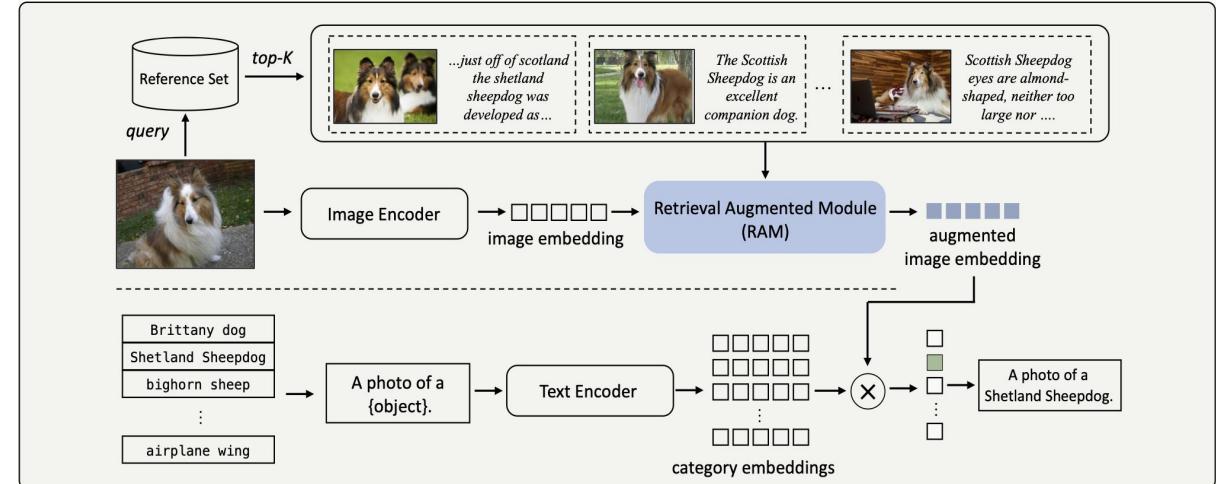
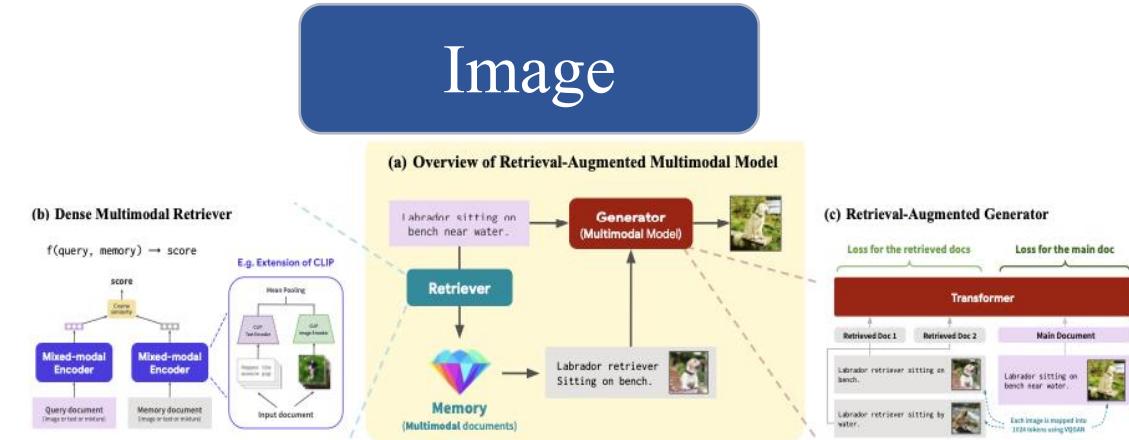
- LLM can be used for **retrieval** (LLM generation replaces retrieval, retrieving from LLM memory), for **generation**, and for **evaluation**. How to further explore the **potential** of LLM in RAG.

Engineering Practice

- How to reduce the **latency** of retrieving ultra-large-scale corpora.
- How to ensure that the content retrieved is not **leaked** by large models

▶ Prospects — Mult-Modality Extension

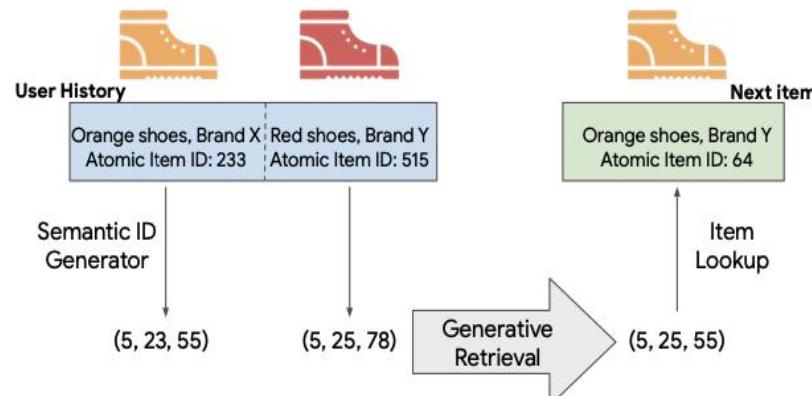
Transferring the concept of RAG from text to other modalities of data



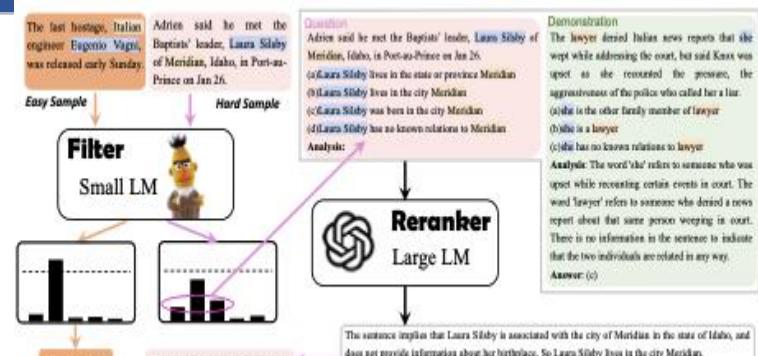
▶ Prospects —— Development of RAG Ecosystem

Further expand the downstream tasks of RAG and improve ecological construction

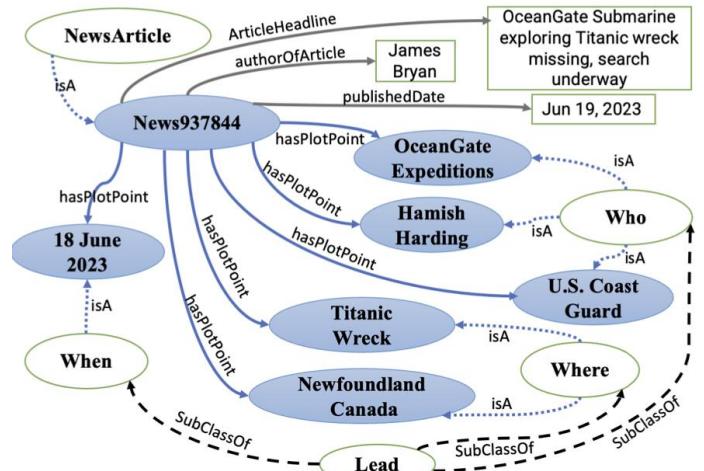
Downstream Task Development and Evaluation



Recommendation System
| TIGER [Rajput et al.,2023]



Information extraction
| Filter- Rerank [Ma et al.,2023]



Report generation
| FABULA [Ranade et al.,2023]

Technology Stack Construction

- **Customized** function, meeting a variety of needs
- **Simplified** use, further reducing the barrier to entry.
- **Specialized** functions, gradually towards production environments.



Personal Knowledge
Assistant Based on RAG



Open-source framework for
production environments