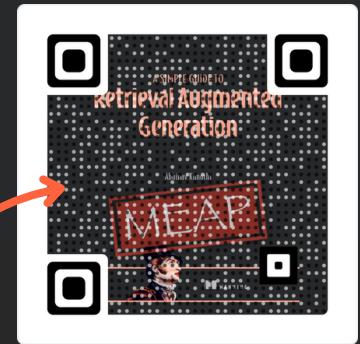


KEY EXCERPTS FROM

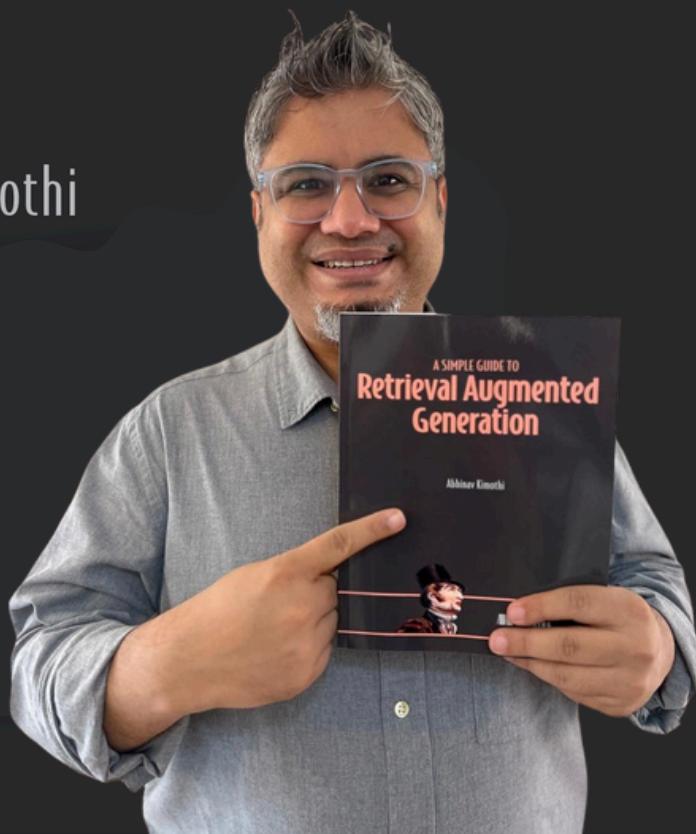


SCAN CODE OR CLICK HERE

A SIMPLE GUIDE TO Retrieval Augmented Generation

Abhinav Kimothi

**50%
OFF**



 MANNING

Table of Contents

Part 1: Foundations

- ① LLMS AND THE NEED FOR RETRIEVAL AUGMENTED GENERATION
- ② RAG SYSTEMS AND THEIR DESIGN

Part 2: Creating RAG Systems

- ③ INDEXING PIPELINE: CREATING A KNOWLEDGE BASE FOR RAG-BASED APPLICATIONS
- ④ GENERATION PIPELINE: GENERATING CONTEXTUAL LLM RESPONSES
- ⑤ RAG EVALUATION: ACCURACY, RELEVANCE, FAITHFULNESS

Part 3: RAG In Production

- ⑥ PROGRESSION OF RAG SYSTEMS: NAÏVE TO ADVANCED, AND MODULAR RAG
- ⑦ EVOLVING RAGOPS STACK: TECHNOLOGIES THAT MAKE RAG POSSIBLE

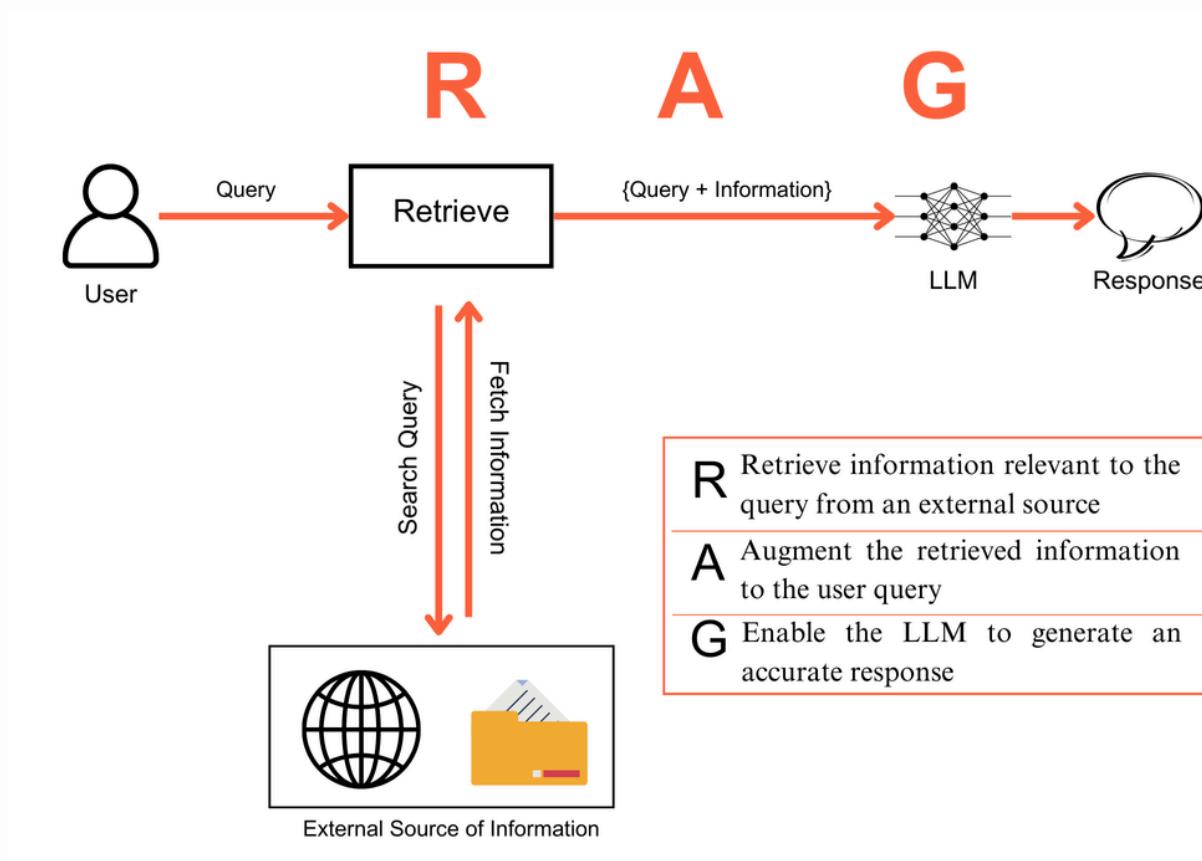
Part 4: Additional Considerations

- ⑧ RAG VARIANTS: MULTIMODAL, AGENTIC, GRAPH AND OTHER RAGS
- ⑨ RAG DEVELOPMENT FRAMEWORK & AREAS OF FURTHER EXPLORATION

Ch 1: LLMs and the Need for Retrieval Augmented Generation

This chapter covers

- Understanding the limits of LLMs & the need for RAG
- Defining Retrieval Augmented Generation (RAG)
- Popular use cases of RAG

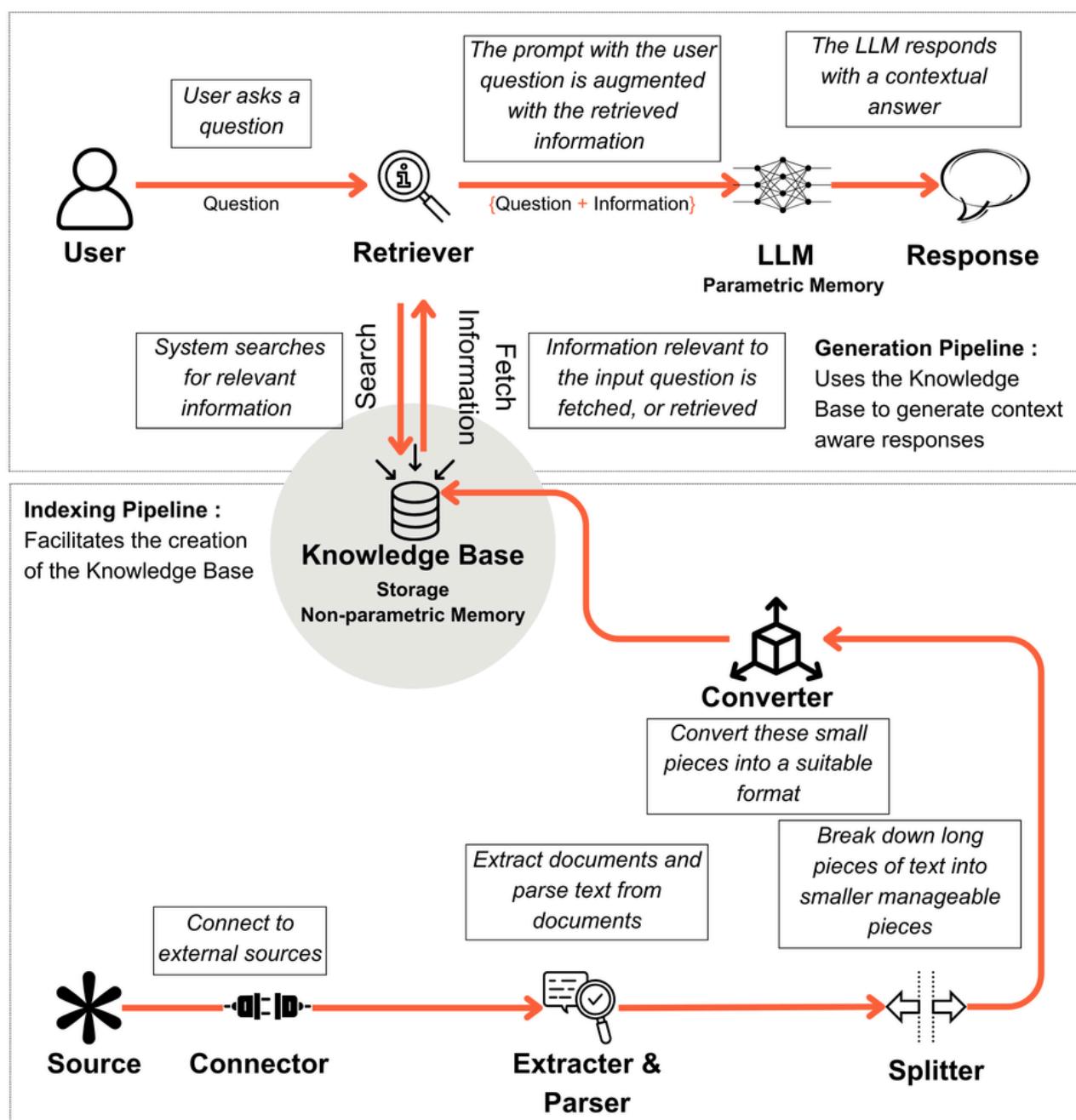


Retrieval Augmented Generation is the methodological approach of enhancing the parametric memory of an LLM by creating access to an explicit non-parametric memory, from which a retriever can fetch relevant information, augment that information to the prompt, pass the prompt to an LLM to enable the LLM to generate a response that is *contextual, reliable, and factually accurate*.

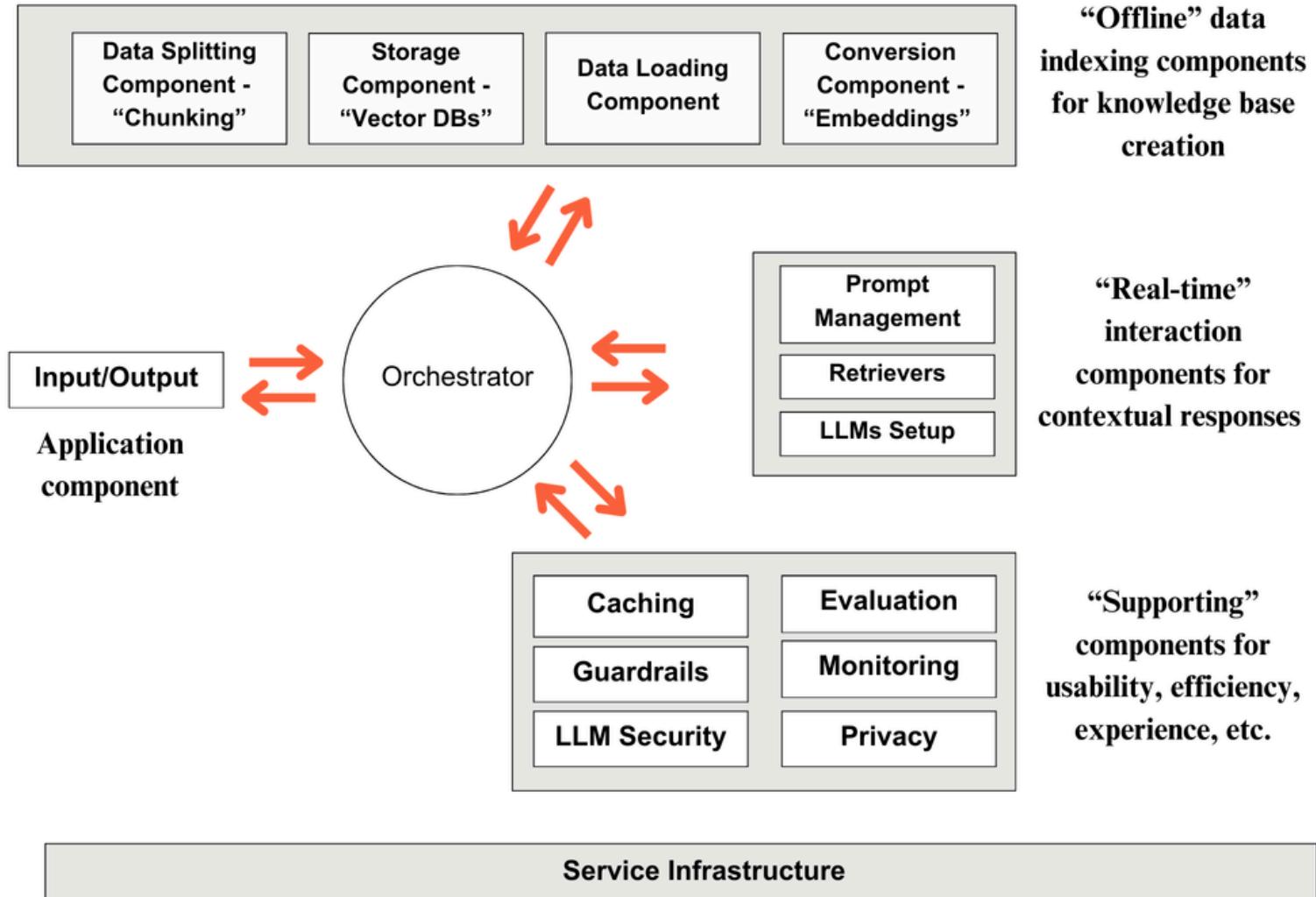
Ch 2 : RAG Systems and Their Design

This chapter covers

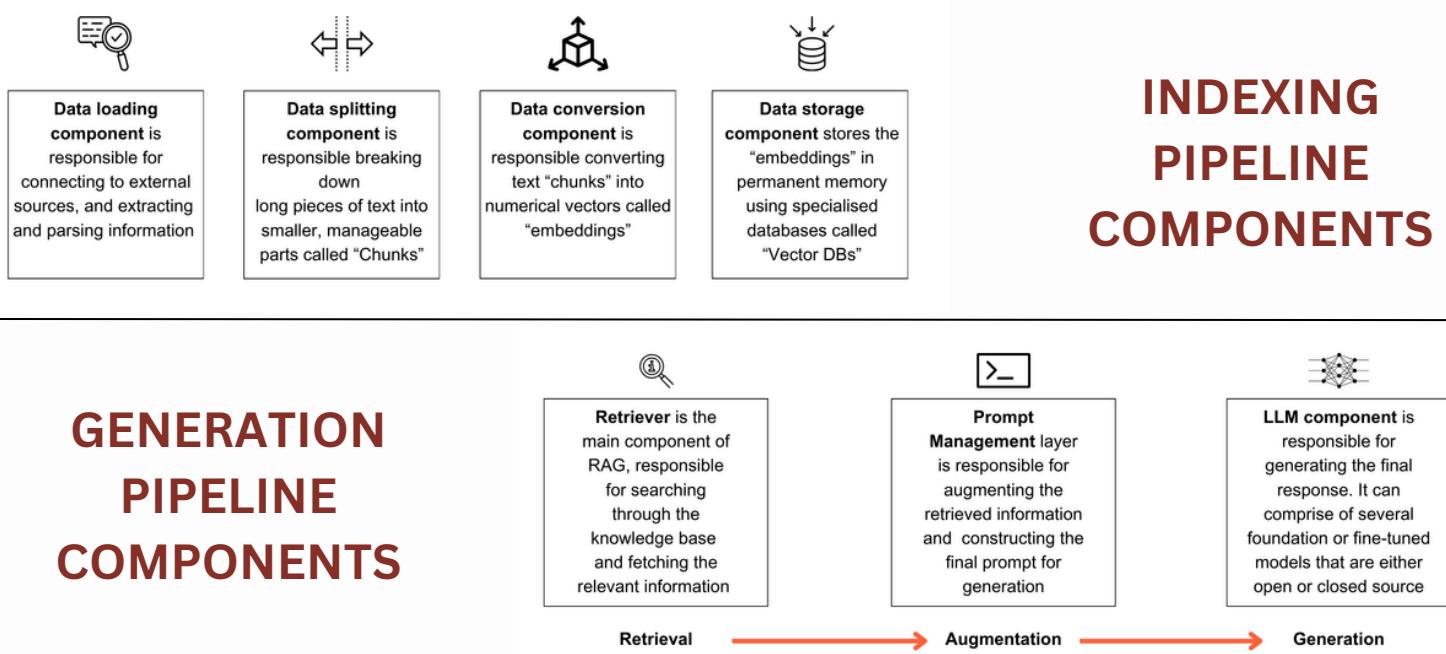
- An exploration of the concept & design of RAG systems
- An Overview of the Indexing Pipeline, Generation Pipeline & their evaluation
- A high-level look at RAG Operations stack



Ch 2 : RAG Systems and Their Design



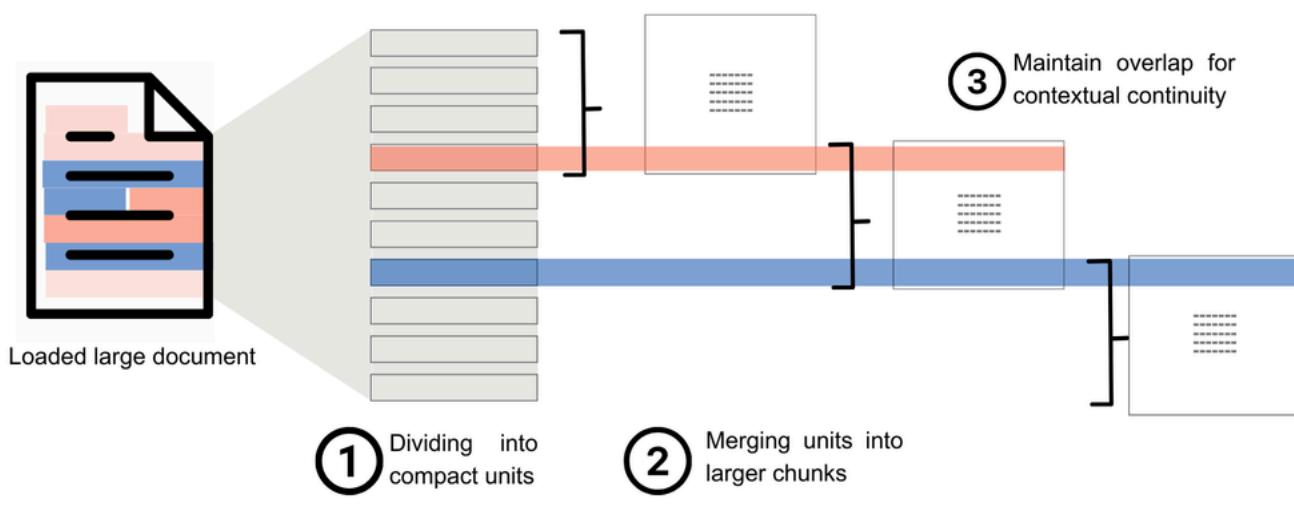
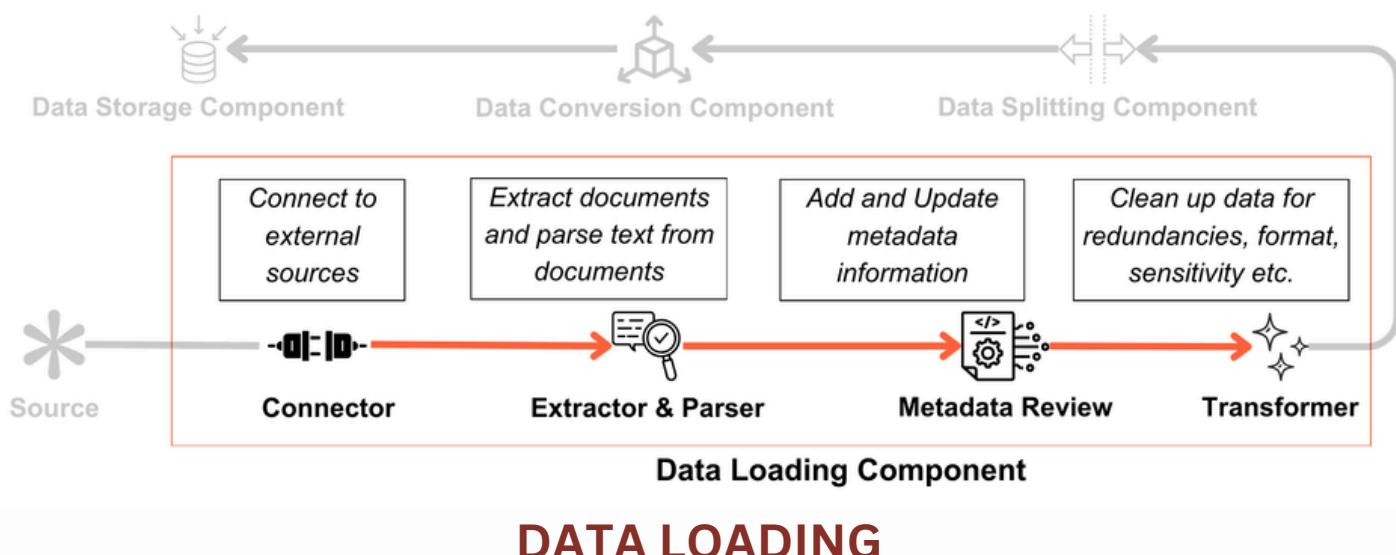
EVOLVING RAGOPS STACK



Ch 3 : Indexing Pipeline: Creating a Knowledge Base for RAG

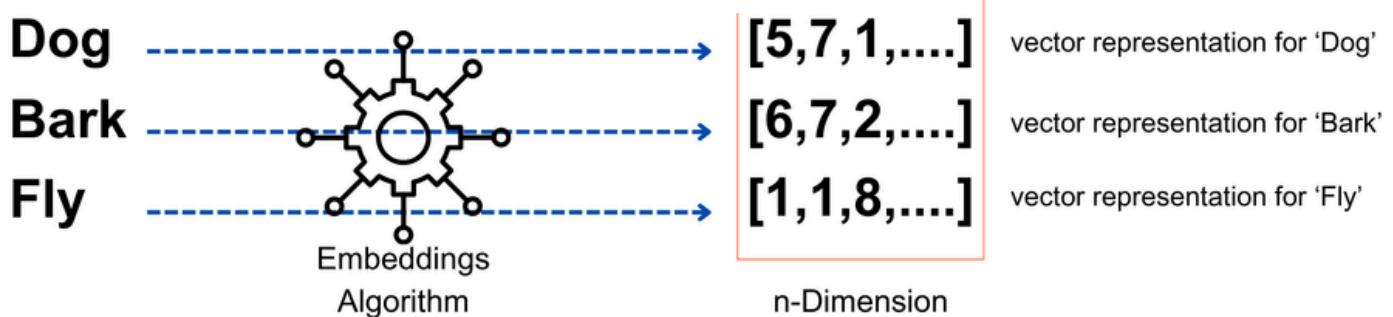
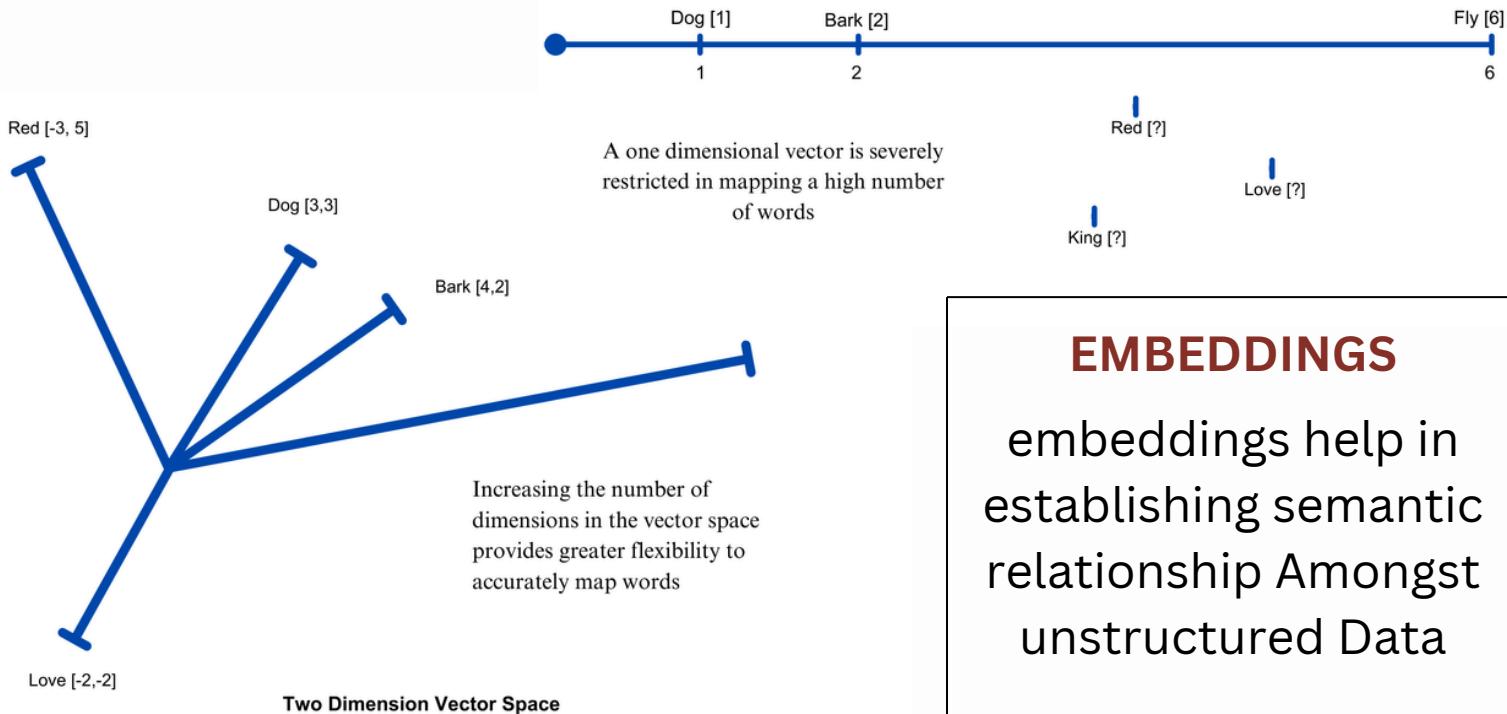
This chapter covers

- The four components of the Indexing Pipeline
 - Data Loading
 - Text Splitting or Chunking
 - Converting Text to Embeddings
 - Storing Embeddings in Vector Databases
- Examples in Python using LangChain.

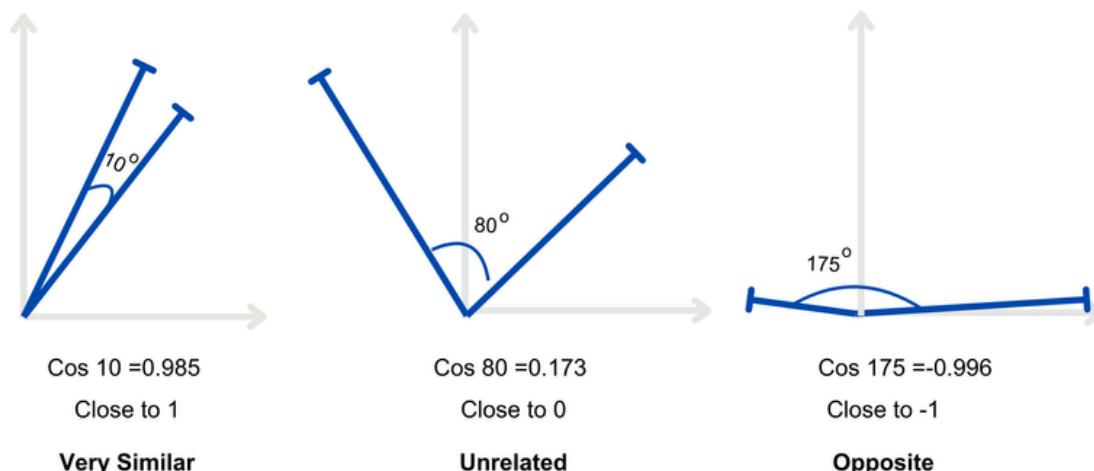


SMALL TO BIG TEXT CHUNKING

Ch 3 : Indexing Pipeline: Creating a Knowledge Base for RAG



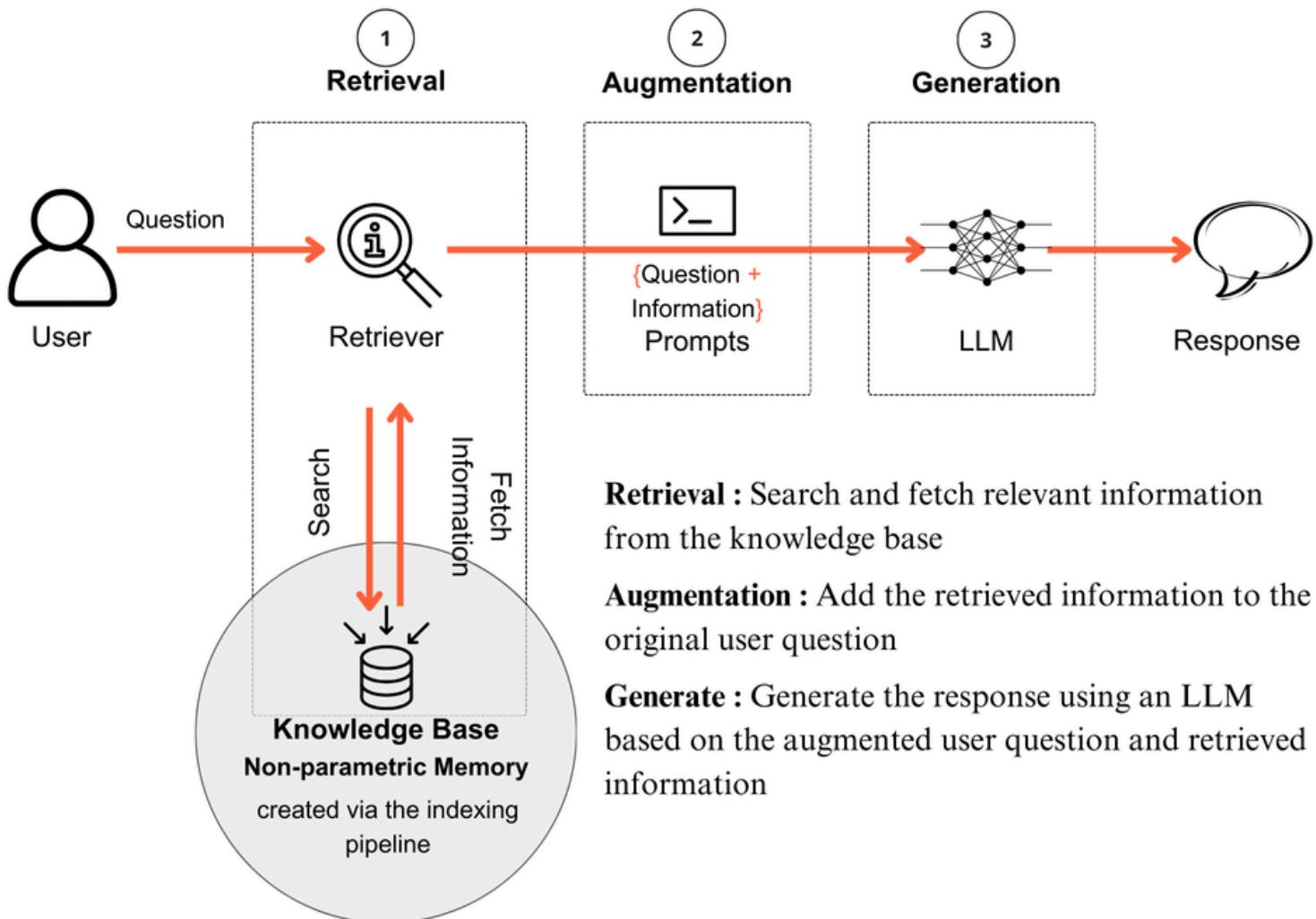
The goal of an embedding model is to convert words (or sentences/paragraphs) into n-dimensional vectors



Ch 4 : Generation Pipeline: Generating Contextual Responses

This chapter covers

- Retrievers and Retrieval Methodologies
- Augmentation with Prompt Engineering Techniques
- Generation using Large Language Models
- Basic implementation of the RAG pipeline in Python



COMPONENTS OF GENERATION PIPELINE

Ch 4 : Generation Pipeline: Generating Contextual Responses

Components of TF-IDF

Term Frequency (TF)

Measures how frequently a term 't' appears in a document 'd'

$$TF(t,d) = \frac{\text{Number of times term } t \text{ appears in document } d}{\text{Total number of terms in document } d}$$

Inverse Document Frequency (IDF)

Measures how important a term 't' is within the entire corpus 'D'

$$IDF(t,D) = \log \left(\frac{\text{Total number of documents 'D'}}{\text{Number of documents containing term 't'}} \right)$$

TF-IDF

Product of TF & IDF

$$TF-IDF(t,d,D) = TF(t,d) \times IDF(t,D)$$

Documents (D)

d1 = Australia won the Cricket World Cup 2023

d2 = India and Australia played in the finals

d3 = Australia won the sixth time & last won in 2015

Search Term

"won"

$$TF("won", d1) = 1/7 = 0.14$$

$$TF("won", d2) = 0/7 = 0$$

$$TF("won", d3) = 2/10 = 0.2$$

$$IDF("won", D) = \log(3/2) = 0.176$$

$$TF-IDF("won", d1, D) = 0.14 \times 0.176 = 0.025$$

$$TF-IDF("won", d2, D) = 0 \times 0.176 = 0$$

$$TF-IDF("won", d3, D) = 0.2 \times 0.176 = 0.035$$

Result - d3 > d1 > d2

Calculating BM25

$$BM25(t,d,D) = IDF(t,D) \times \frac{TF(t,d) \times (k+1)}{TF(t,d) + (k) \times (1-b + b \times \frac{|d|}{avgdl})}$$

- $TF(t,d)$ is the term frequency of term 't' in document 'd'
- $IDF(t,D)$ is the inverse document frequency of term in the corpus
- $|d|$ is the length of the document
- $avgdl$ is the average document length in the entire corpus.
- k and b are free parameters

Documents (D)

d1 = Australia won the Cricket World Cup 2023

d2 = India and Australia played in the finals

d3 = Australia won the sixth time & last won in 2015

$$BM25("won", d1, D) = 0.193$$

$$BM25("won", d2, D) = 0$$

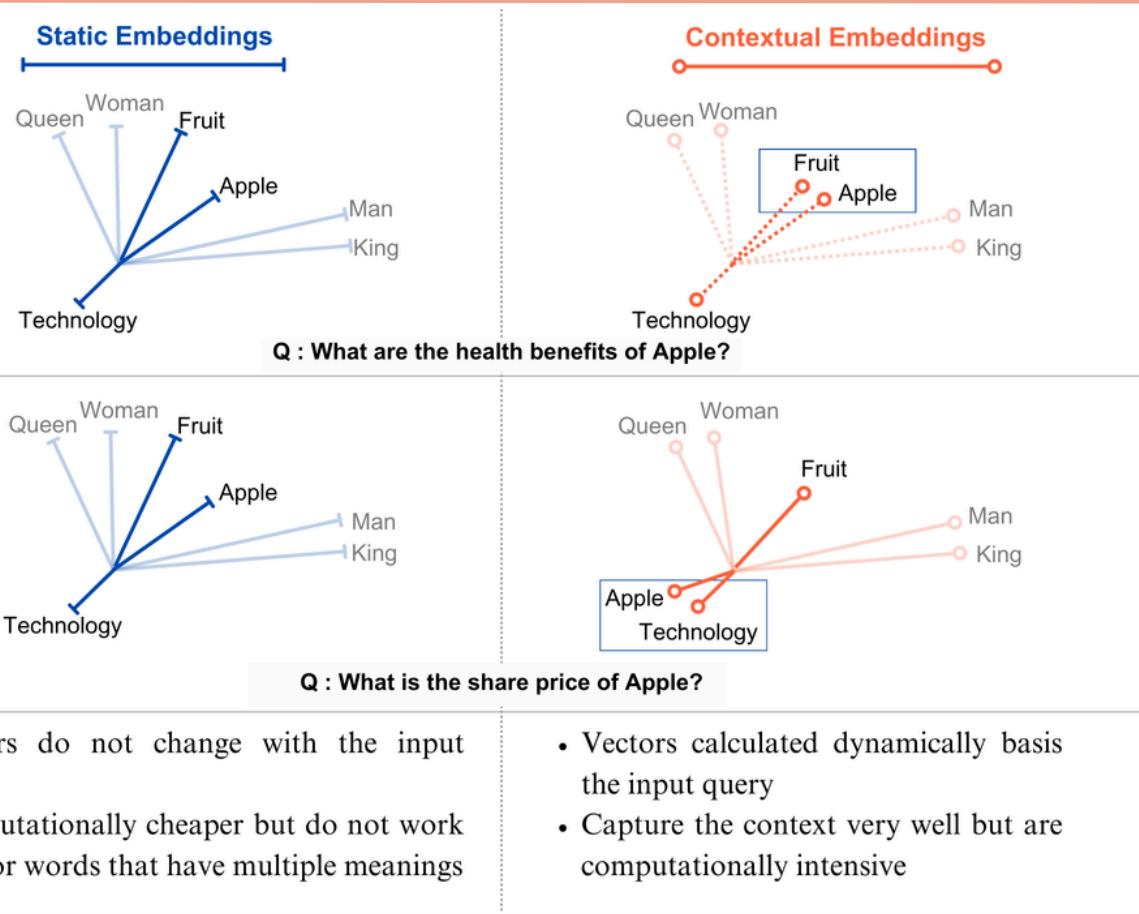
$$BM25("won", d3, D) = 0.168$$

Result - d1 > d3 > d2

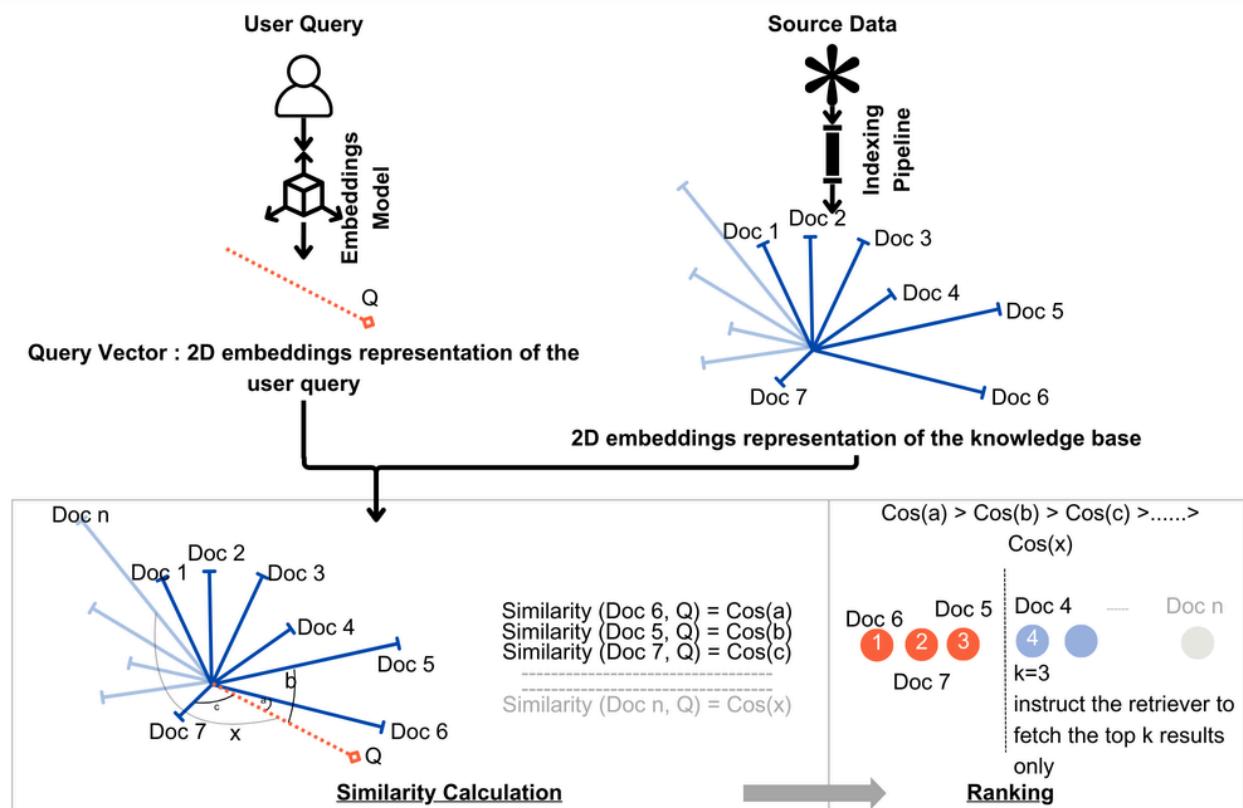
TF-IDF AND BM25 ARE EXAMPLES OF SPARSE RETRIEVERS

Ch 4 : Generation Pipeline: Generating Contextual Responses

STATIC VS CONTEXTUAL EMBEDDINGS



COSINE SIMILARITY BASED DOCUMENT RETRIEVAL

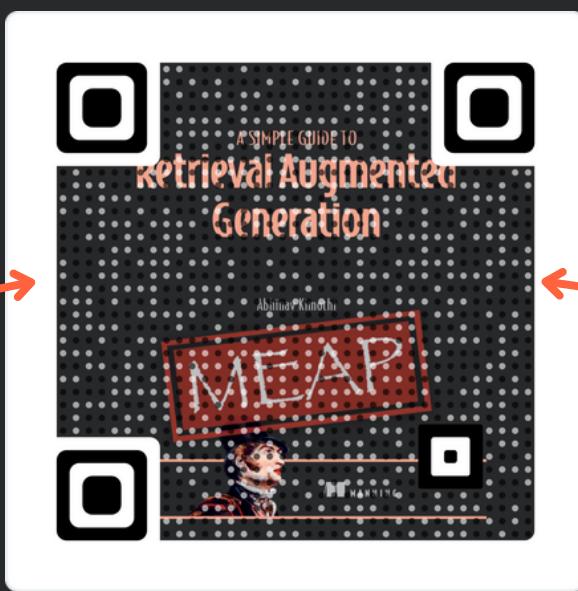
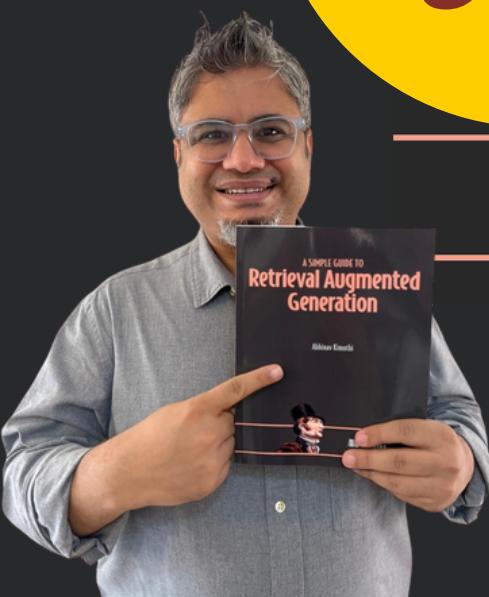


THESE EXCERPTS ARE FROM A SIMPLE GUIDE TO RETRIEVAL AUGMENTED GENERATION

A SIMPLE GUIDE TO
**Retrieval Augmented
Generation**

**50%
OFF**

1400+
Copies Sold



**SCAN CODE OR CLICK HERE
TO GET AN EARLY ACCESS COPY**

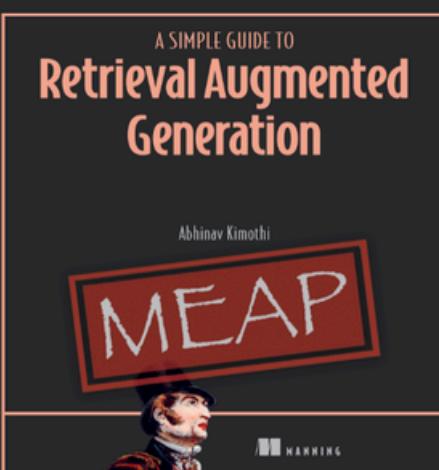
INTERESTED IN CODING RAG PIPELINES?

THE SOURCE CODE OF

A SIMPLE GUIDE TO

RETRIEVAL AUGMENTED GENERATION

IS NOW AVAILABLE FOR FREE PUBLIC ACCESS



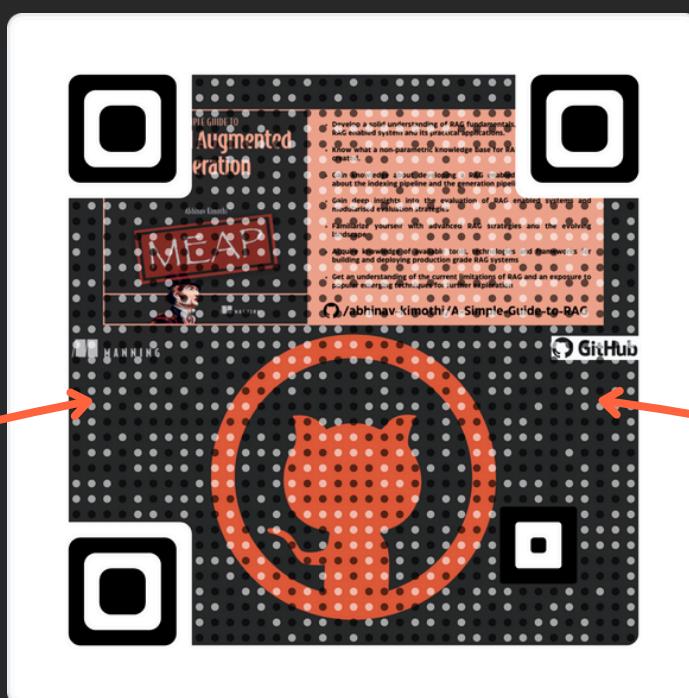
- Develop a solid understanding of RAG fundamentals, the components of a RAG enabled system and its practical applications.
- Know what a non-parametric knowledge base for RAG means and how is it created.
- Gain knowledge about developing a RAG enabled system with details about the indexing pipeline and the generation pipeline.
- Gain deep insights into the evaluation of RAG enabling modularised evaluation strategies
- Familiarize yourself with advanced RAG strategies and the existing landscape
- Acquire knowledge of available tools, technologies and frameworks for building and deploying production grade RAG systems
- Get an understanding of the current limitations of RAG and an exposure to popular emerging techniques for further exploration

[/abhinav-kimothi/A-Simple-Guide-to-RAG](https://github.com/abhinav-kimothi/A-Simple-Guide-to-RAG)

180+ stars

MANNING

GitHub

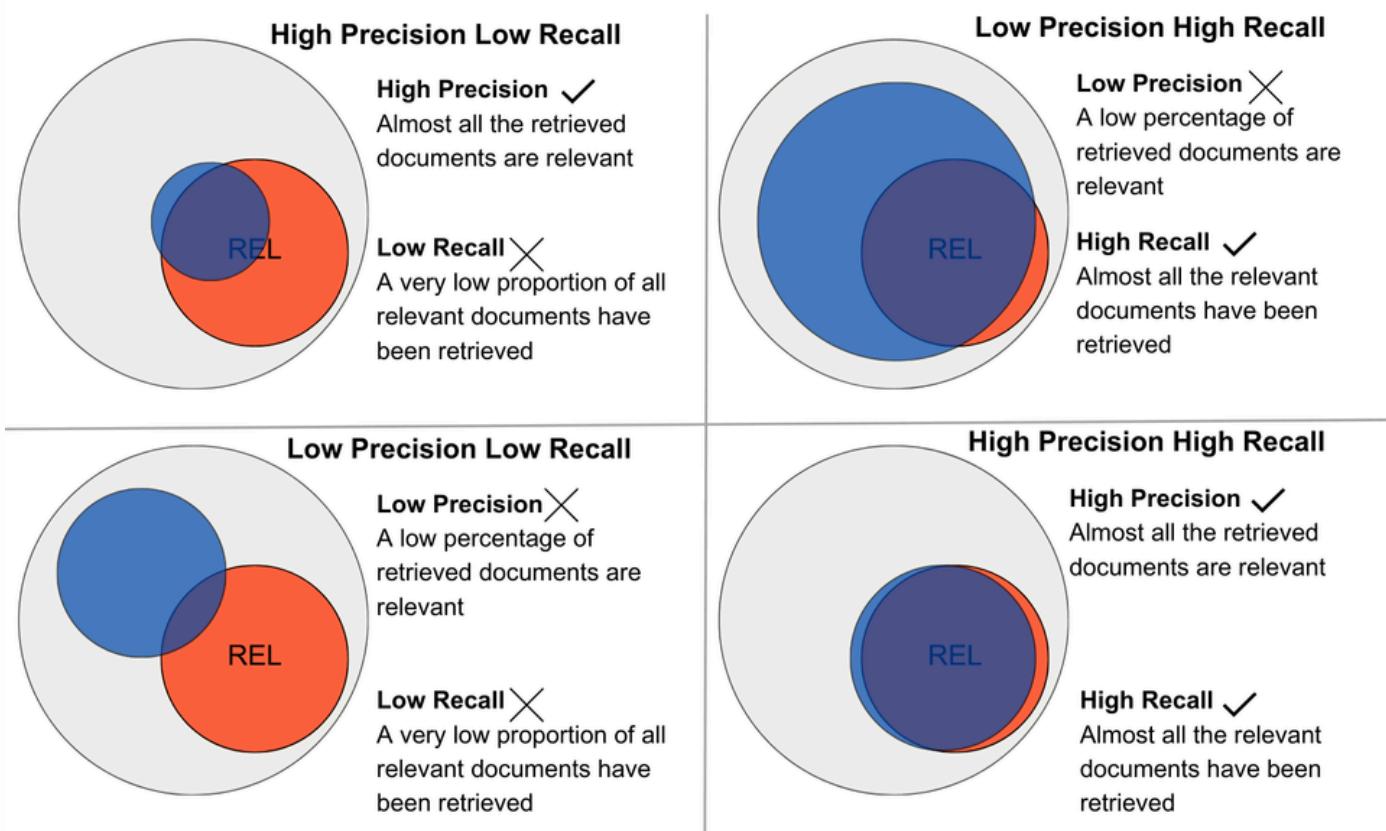
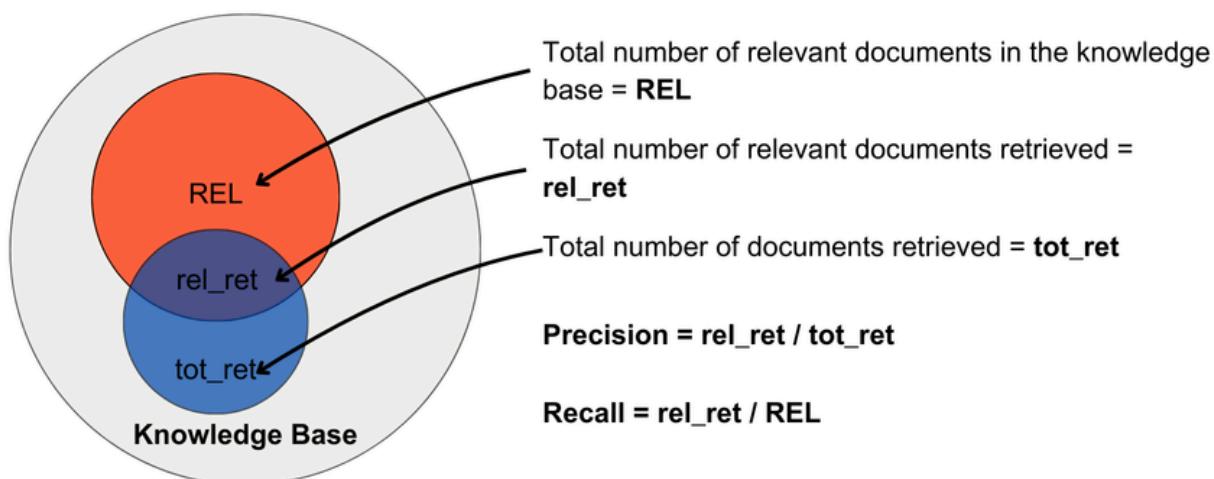


SCAN CODE OR CLICK HERE
TO VIEW SOURCE CODE

Ch 5 : RAG Evaluation: Accuracy, Relevance, Faithfulness

This chapter covers

- The need and requirements for evaluating RAG pipelines
- Metrics, Frameworks, and Benchmarks for RAG Evaluation
- Current limitations and future course of RAG evaluation



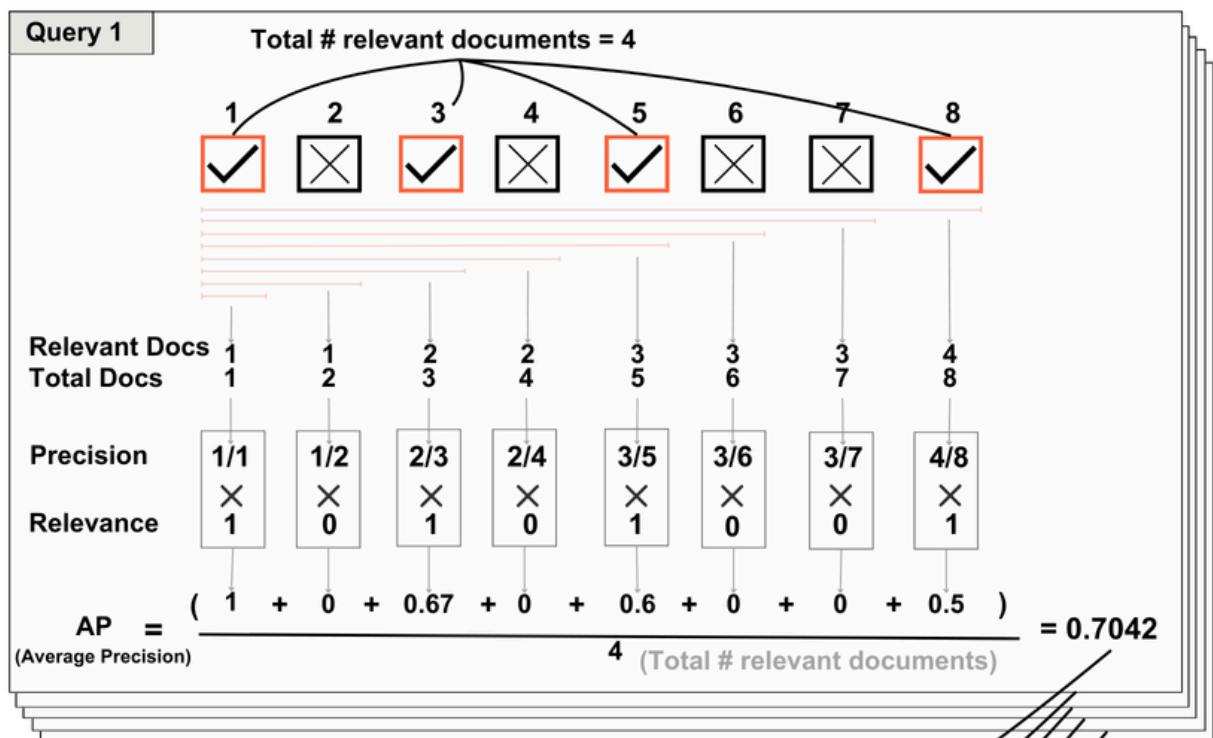
PRECISION AND RECALL SCENARIOS IN DOCUMENT RETRIEVAL

Ch 5 : RAG Evaluation: Accuracy, Relevance, Faithfulness

	1	2	3	4	5	Rank of 1st relevant	Reciprocal
Query 1:	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	3	1/3
Query 2:	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	1	1
	A relevant result on rank 1 results in perfect reciprocal rank						
Query 3:	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	2	1/2
	Considers only the first relevant result						
Query 4:	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	3	1/3
	Doesn't account for number of relevant results						
MRR =	$\frac{1/3 + 1 + 1/2 + 1/3}{4} = 13/24 = 0.54$						

MEAN RECIPROCAL RANK

MEAN AVERAGE PRECISION



Ch 5 : RAG Evaluation: Accuracy, Relevance, Faithfulness

Query : Who won the 2023 ODI Cricket World Cup and when?

Context 1 : High Context Relevance

The 2023 Cricket World Cup, concluded on 19 November 2023, with Australia winning the tournament. The tournament took place in ten different stadiums, in ten cities across the country.

Total sentences = 2

Relevant sentences = 1

Context Relevance = 0.5 or 50%

Context 2 : Low Context Relevance

The 2023 Cricket World Cup was the 13th edition of the Cricket World Cup. It was the first Cricket World Cup which India hosted solely. The tournament took place in ten different stadiums. In the first semi-final India beat New Zealand, and in the second semi-final Australia beat South Africa.

Total sentences = 4

Relevant sentences = 0

Context Relevance = 0

CONTEXT RELEVANCE

Query : Who won the 2023 ODI Cricket World Cup and when?

Context : The 2023 ODI Cricket World Cup concluded on 19 November 2023, with Australia winning the tournament.

Response 1 : High Faithfulness

[Australia] won on [19 November 2023]

Number of claims generated = 2

Number of claims in context = 2

Answer Faithfulness = 1 or 100%

Response 2 : Low Faithfulness

[Australia] won on [15 October 2023] by [defeating India]

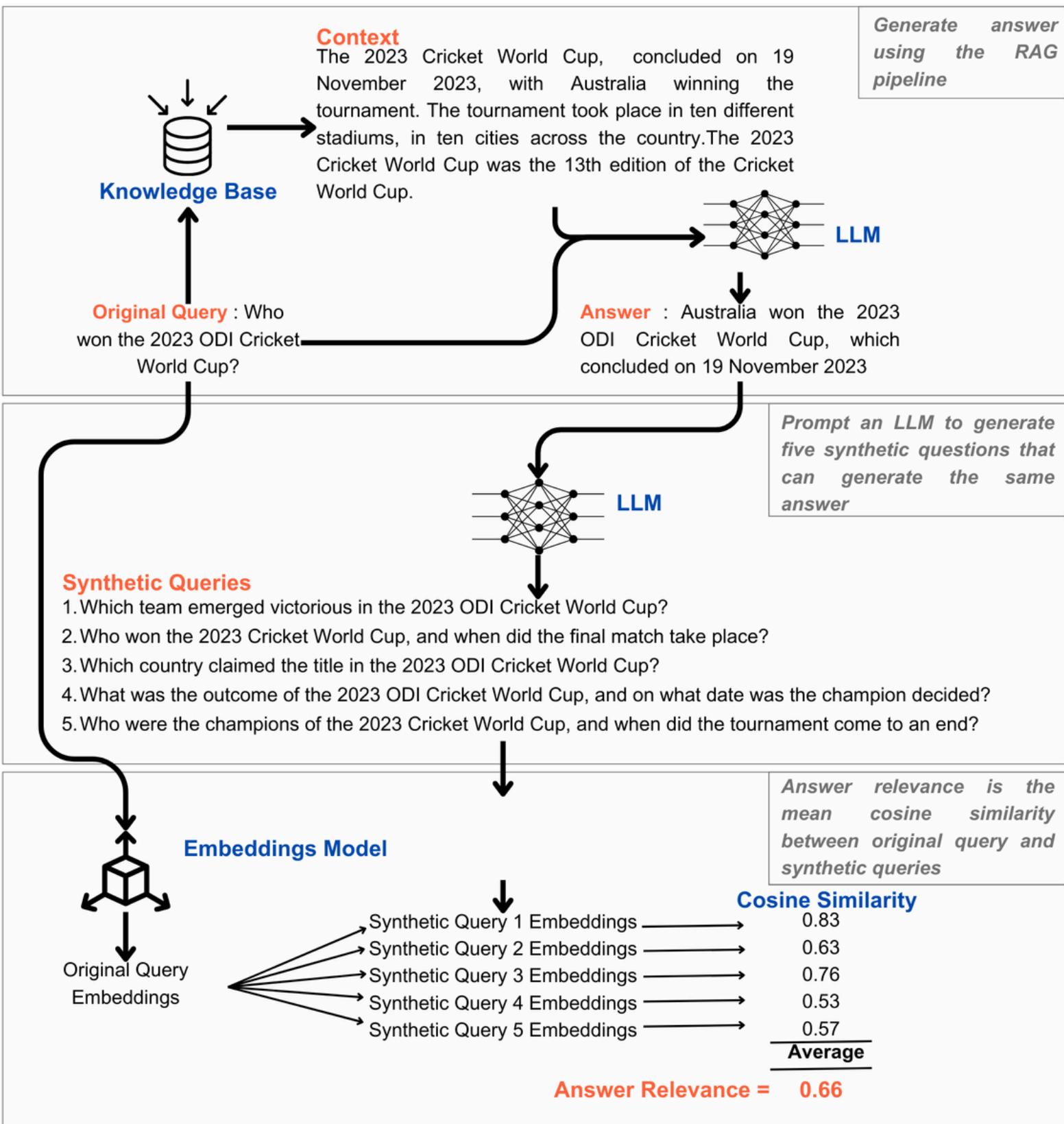
Number of claims generated = 3

Number of claims in context = 1

Answer Faithfulness = 0.33 or 33%

FAITHFULNESS OR GROUNDEDNESS

Ch 5 : RAG Evaluation: Accuracy, Relevance, Faithfulness

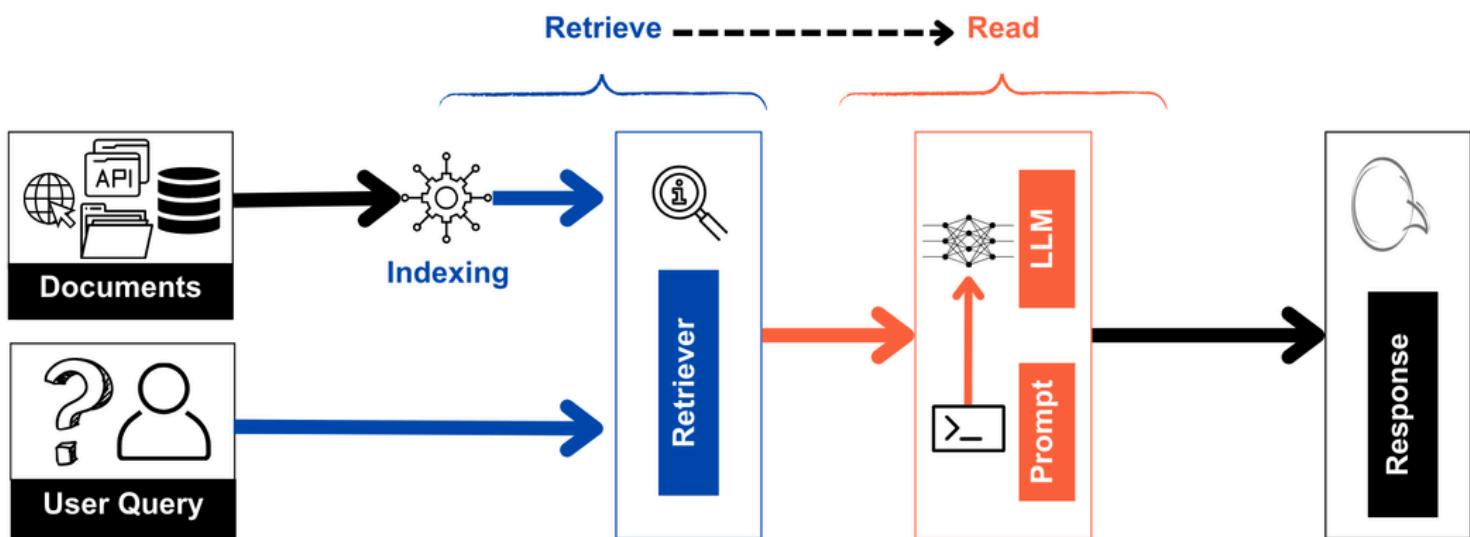


ANSWER RELEVANCE

Ch 6 : RAG Systemd Progression: Naïve, Advanced, & Modular RAG

This chapter covers

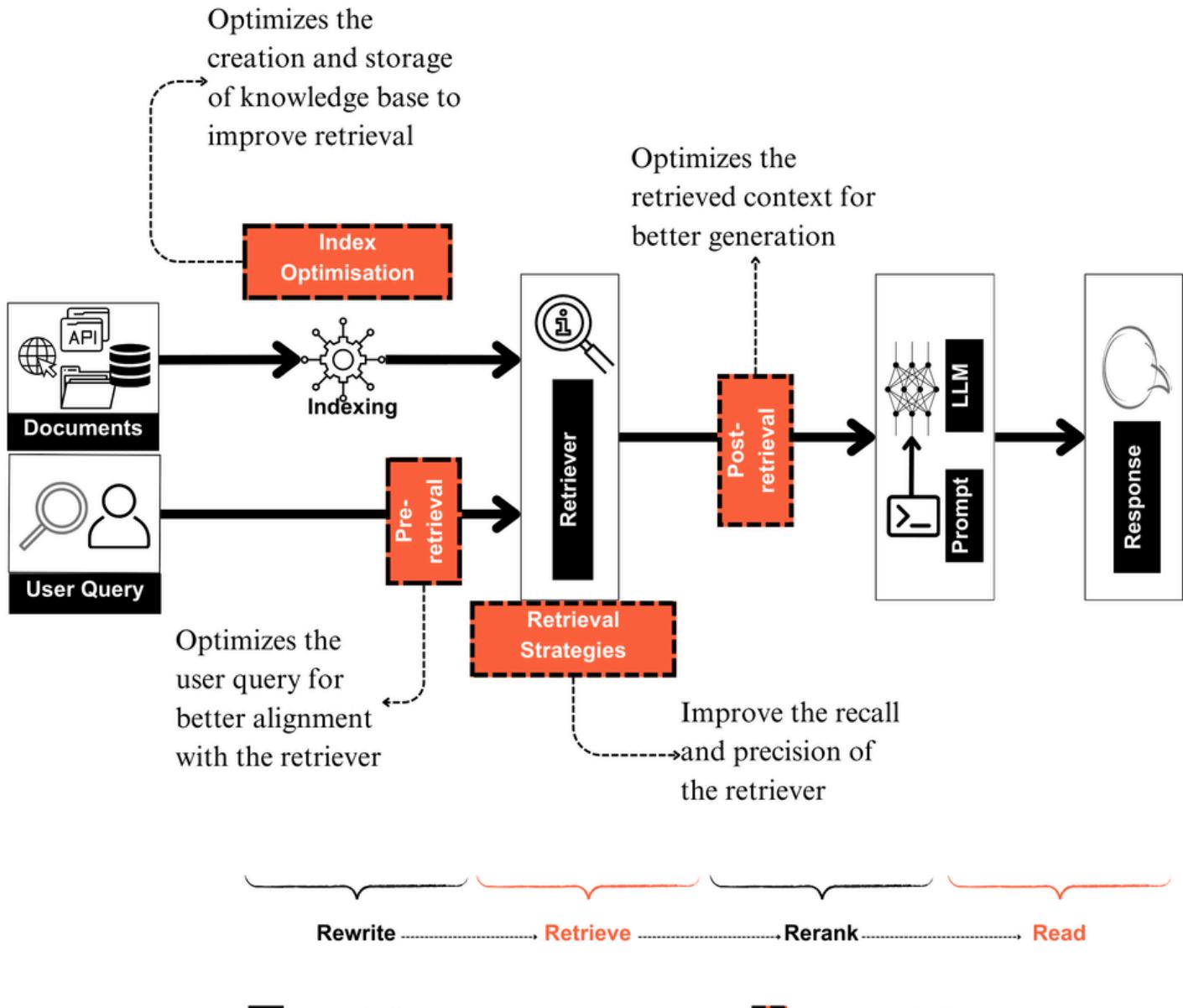
- Limitations of the Naïve RAG approach
- Advanced RAG strategies and techniques
- Modular patterns in RAG



LIMITATIONS OF NAÏVE RAG

The Naïve RAG approach falls in a “retrieve then read” framework - which means that there’s a retriever that is retrieving information and then there’s an LLM that is reading this information to generate the results. This process is marred with drawbacks at each of the three stages of retrieval, augmentation and generation

Ch 6 : RAG Systemd Progression: Naïve, Advanced, & Modular RAG



ADVANCED RAG

With techniques employed at the three stages, the Advanced RAG process follows a ‘Rewrite then Retrieve then Re-rank then Read’ frameworks. Two additional components of Rewrite and Re-rank are added, and the Retrieve component is enhanced in comparison with Naïve RAG.

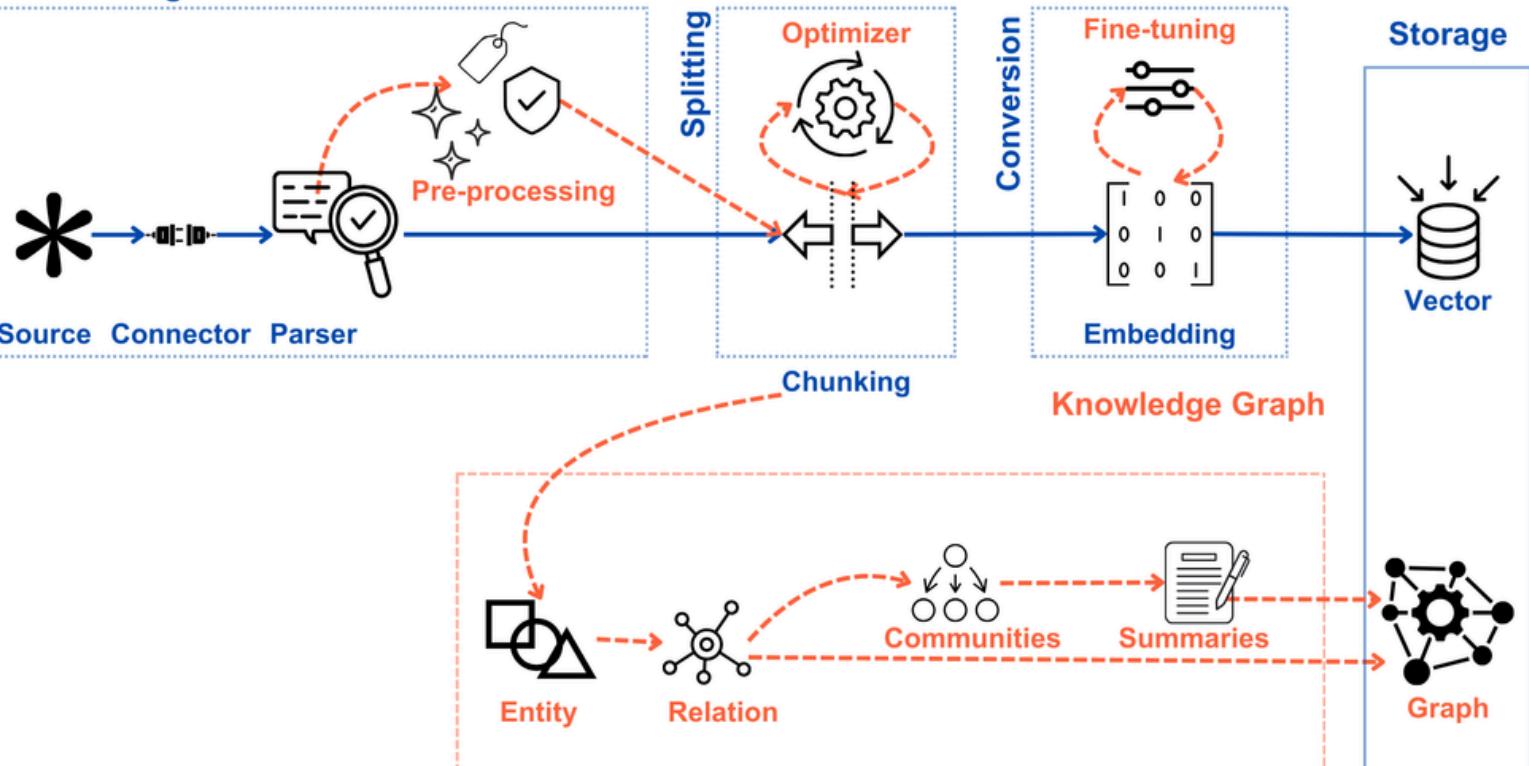
Ch 6 : RAG Systemd Progression: Naïve, Advanced, & Modular RAG

Pre-processing steps like cleaning, metadata enhancement while loading data enhance the searchability

Chunk Size optimization and enriched context enhances retrieval and contextual generation

Embeddings model can be fine-tuned for domain adaptability

Data Loading



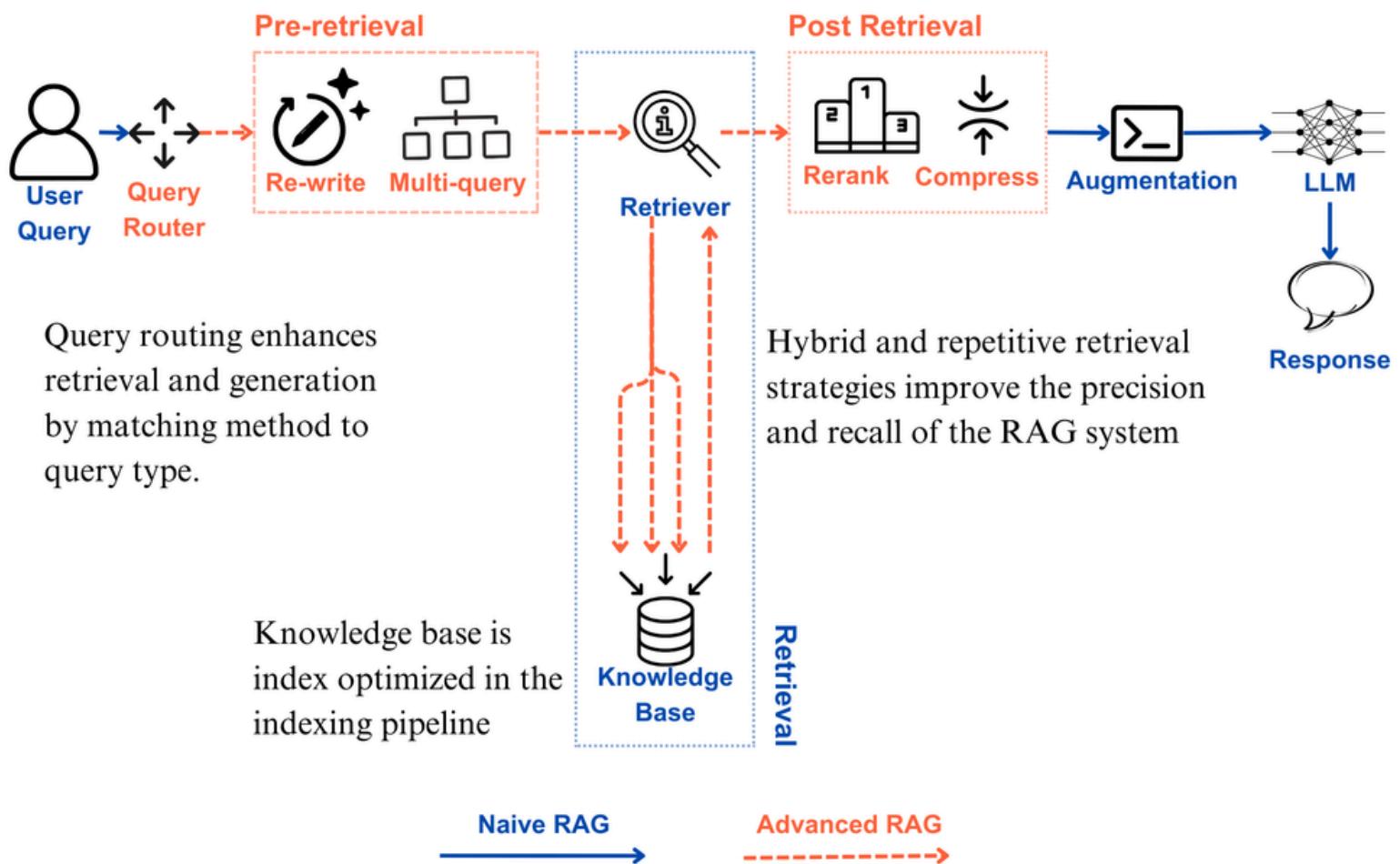
Knowledge Graph Indexing results in deeper context by establishing entity relationships and creating higher order community sub-graph summaries

INDEXING PIPELINE MODIFIED FOR ADVANCED RAG

Ch 6 : RAG Systemd Progression: Naïve, Advanced, & Modular RAG

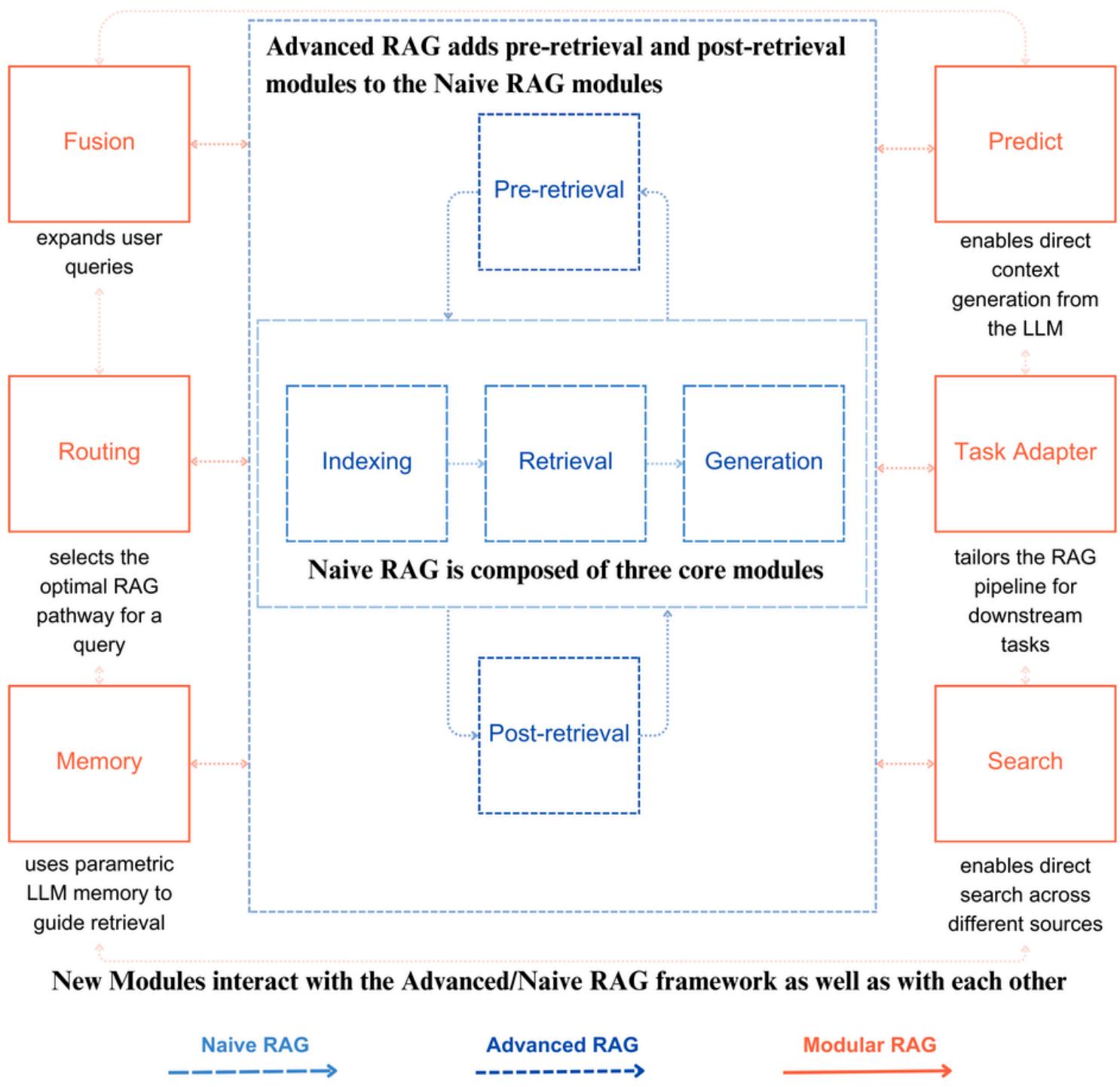
At the pre-retrieval stage, interventions optimize the input user query in a manner that makes it better suited for the retrieval tasks

Post-retrieval techniques optimize the retrieved context for better alignment with the model which results in more coherent and contextual responses



GENERATION PIPELINE MODIFIED FOR ADVANCED RAG

Ch 6 : RAG Systemd Progression: Naïve, Advanced, & Modular RAG



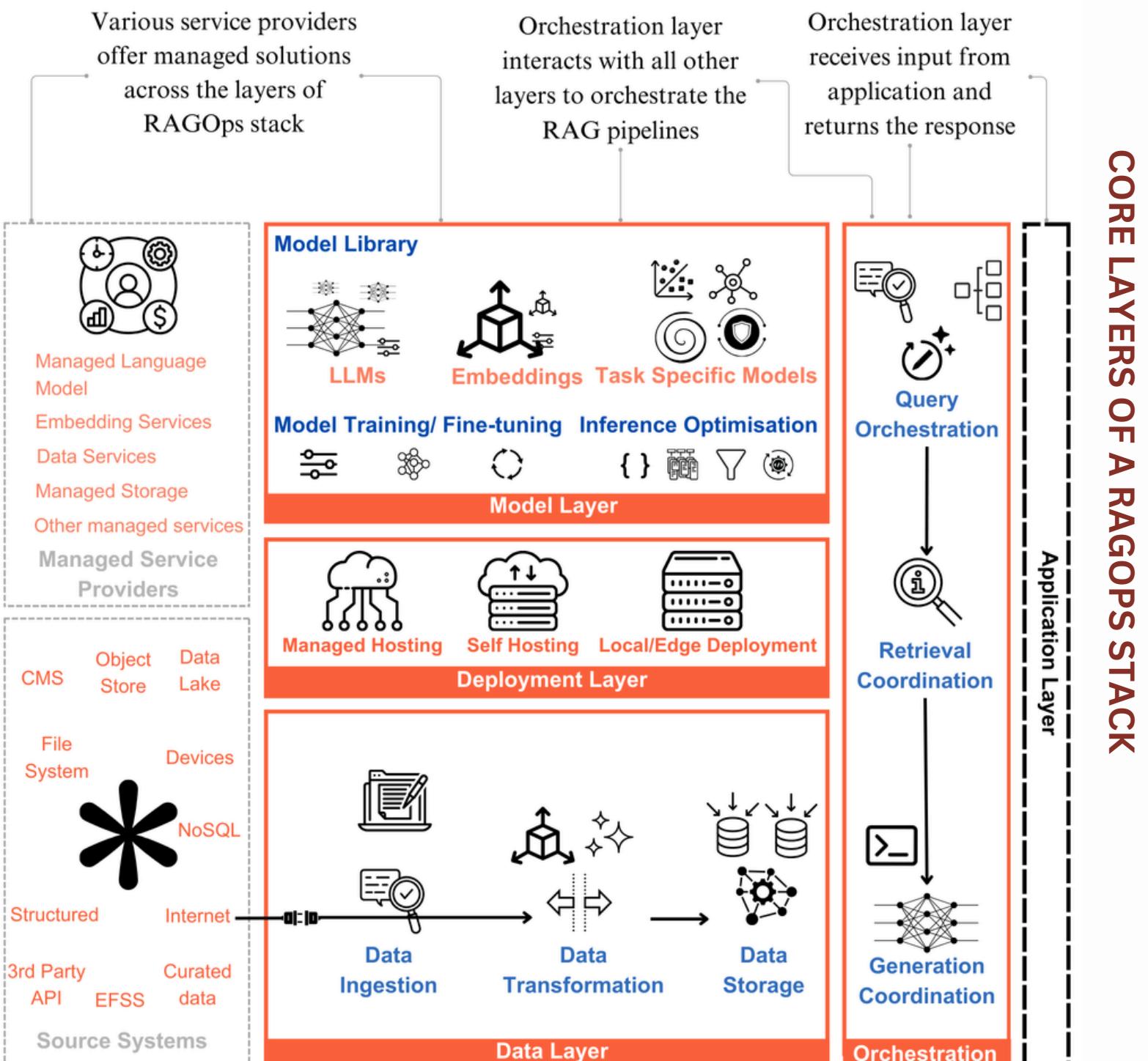
GENERATION PIPELINE MODIFIED FOR ADVANCED RAG

Naïve, Advanced, and Modular approaches to RAG are progressive. Naïve RAG is a sub-component of Advanced RAG which is a sub-component of Modular RAG.

Ch 7 : Evolving RAGOps Stack: Making RAG possible

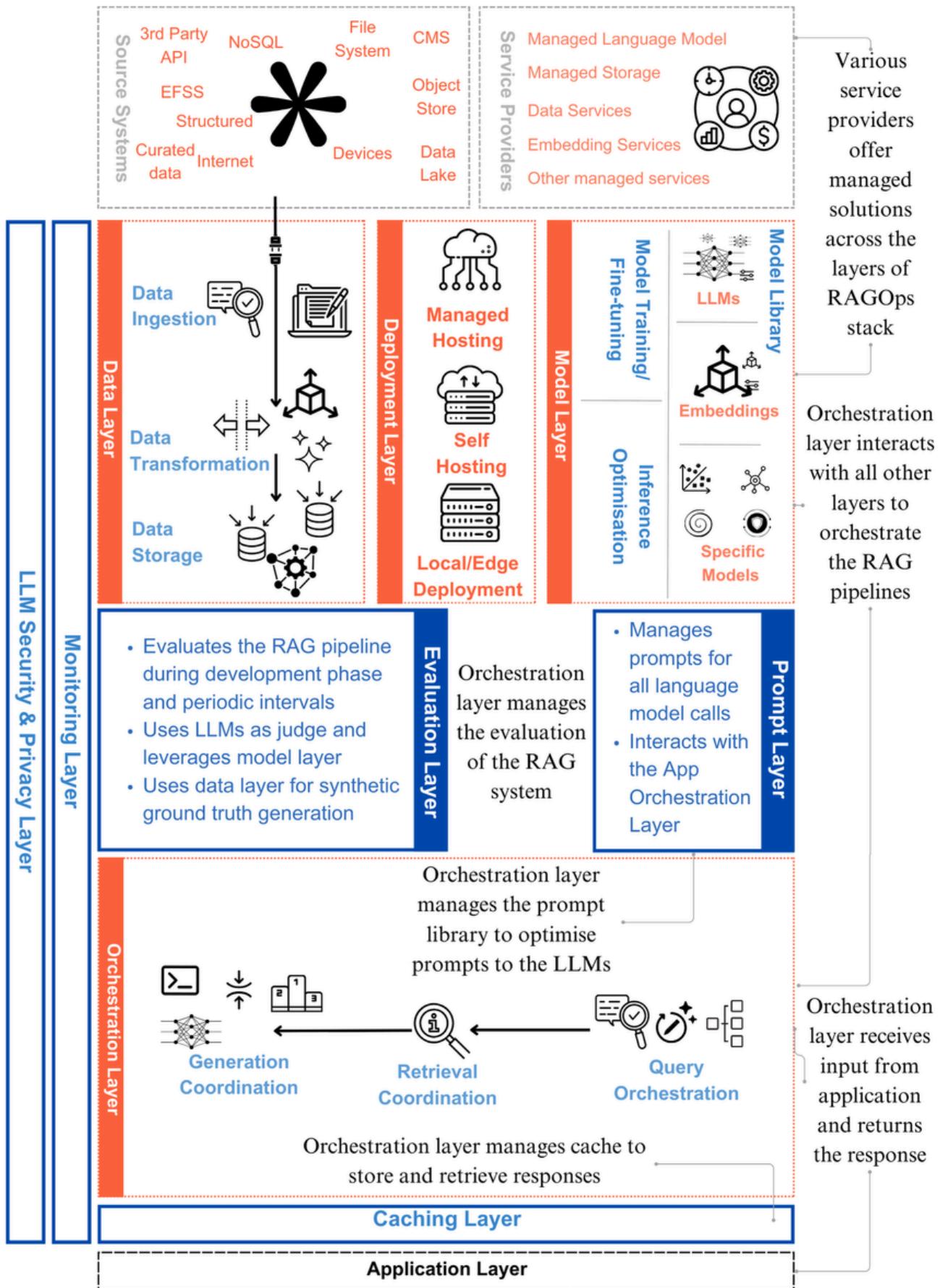
This chapter covers

- The design of RAG systems
- Available tools and technologies that enable a RAG system
- Production best practices for RAG systems



Ch 7 : Evolving RAGOps Stack: Making RAG possible

CORE & ESSENTIAL LAYERS OF A RAGOPS STACK

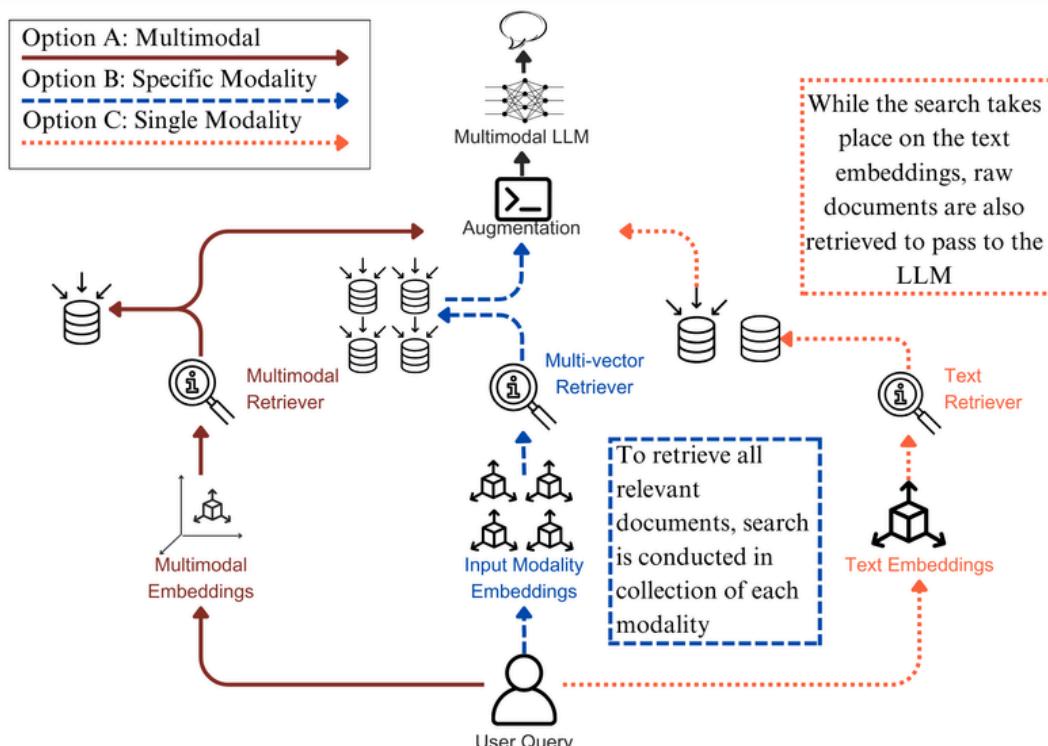
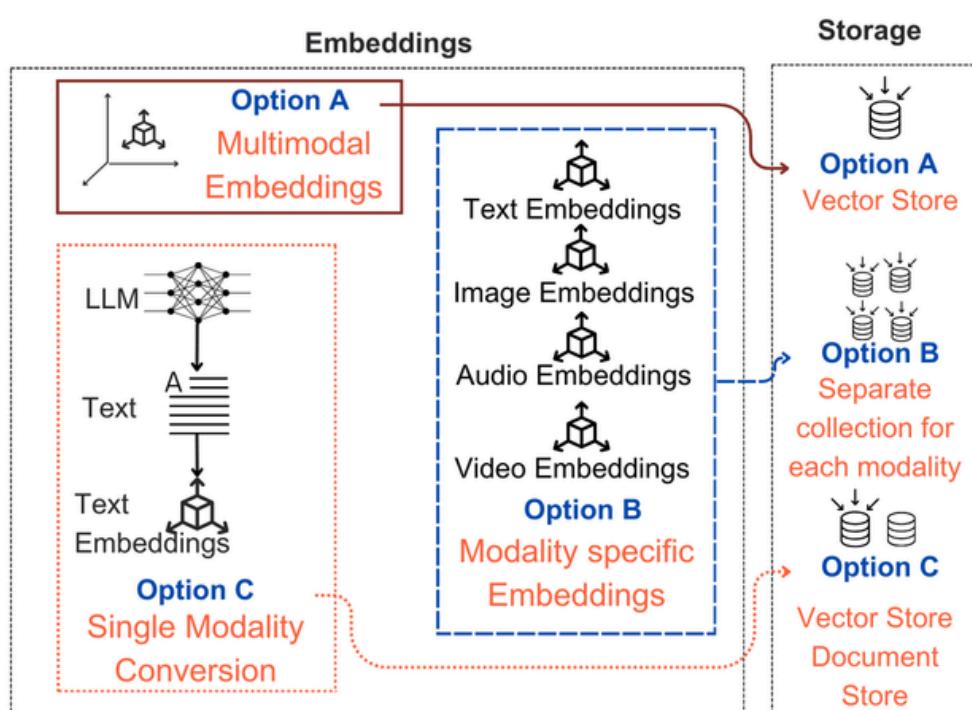
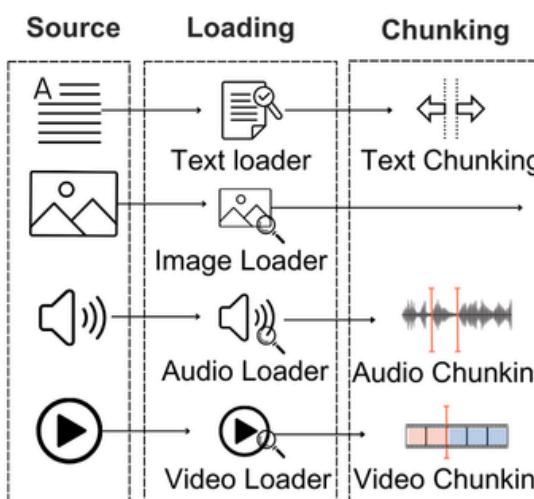


Ch 8 : RAG variants: Multimodal, agentic, graph and other RAGs

This chapter covers

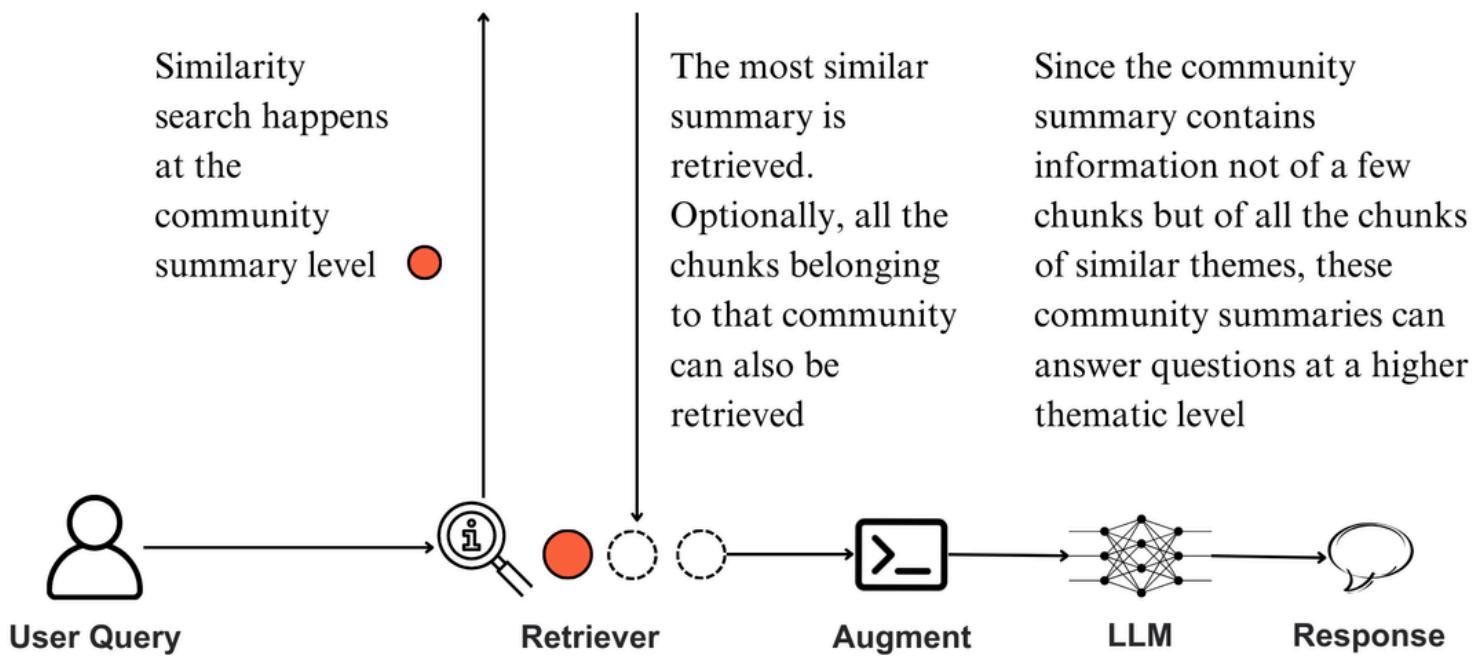
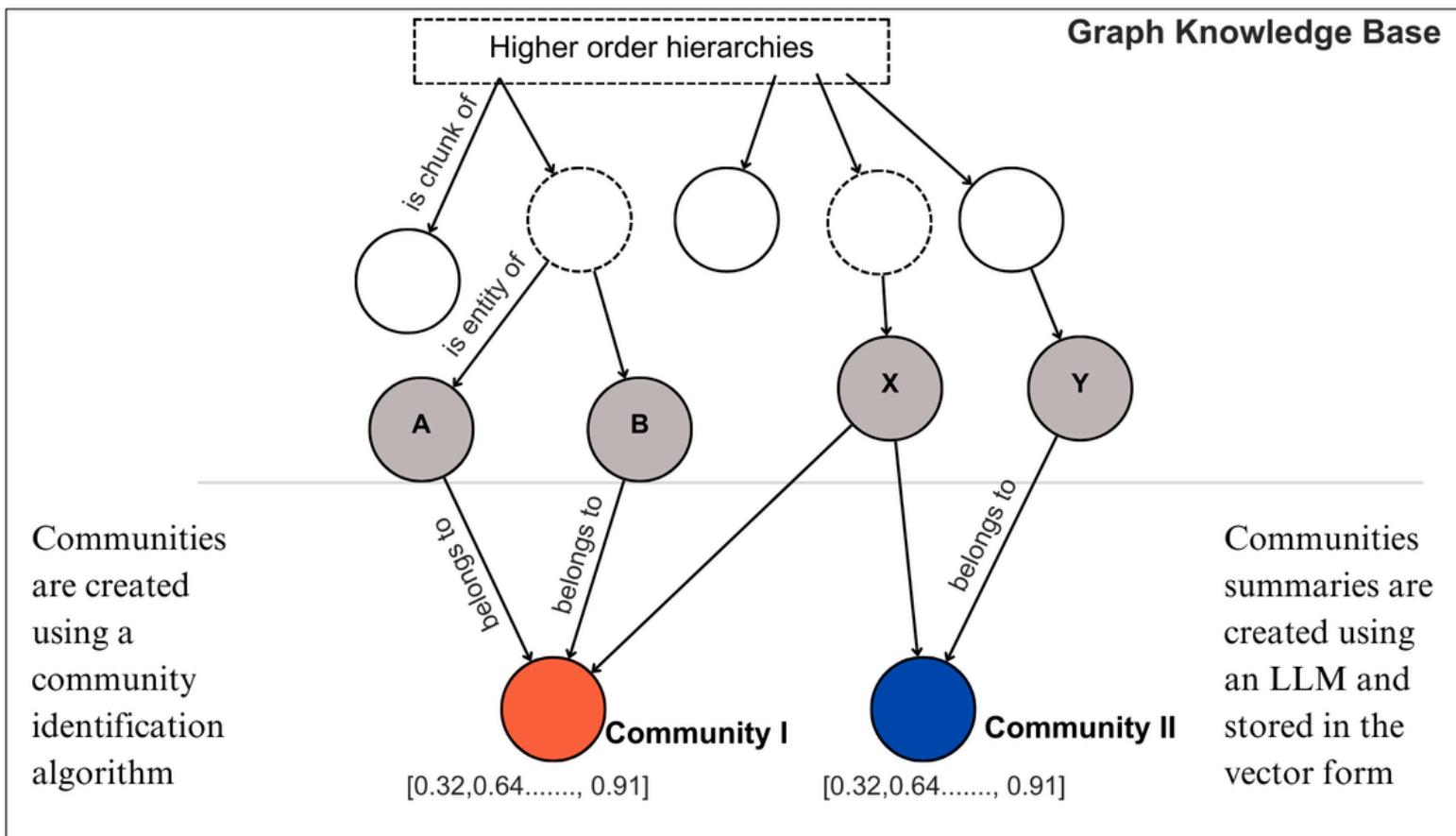
- Introduction to RAG Variants
- Knowledge Graph RAG, Multimodal RAG, Agentic RAG
- Other RAG Variants

Loading and Chunking approach remains largely similar for each of the multimodal RAG options



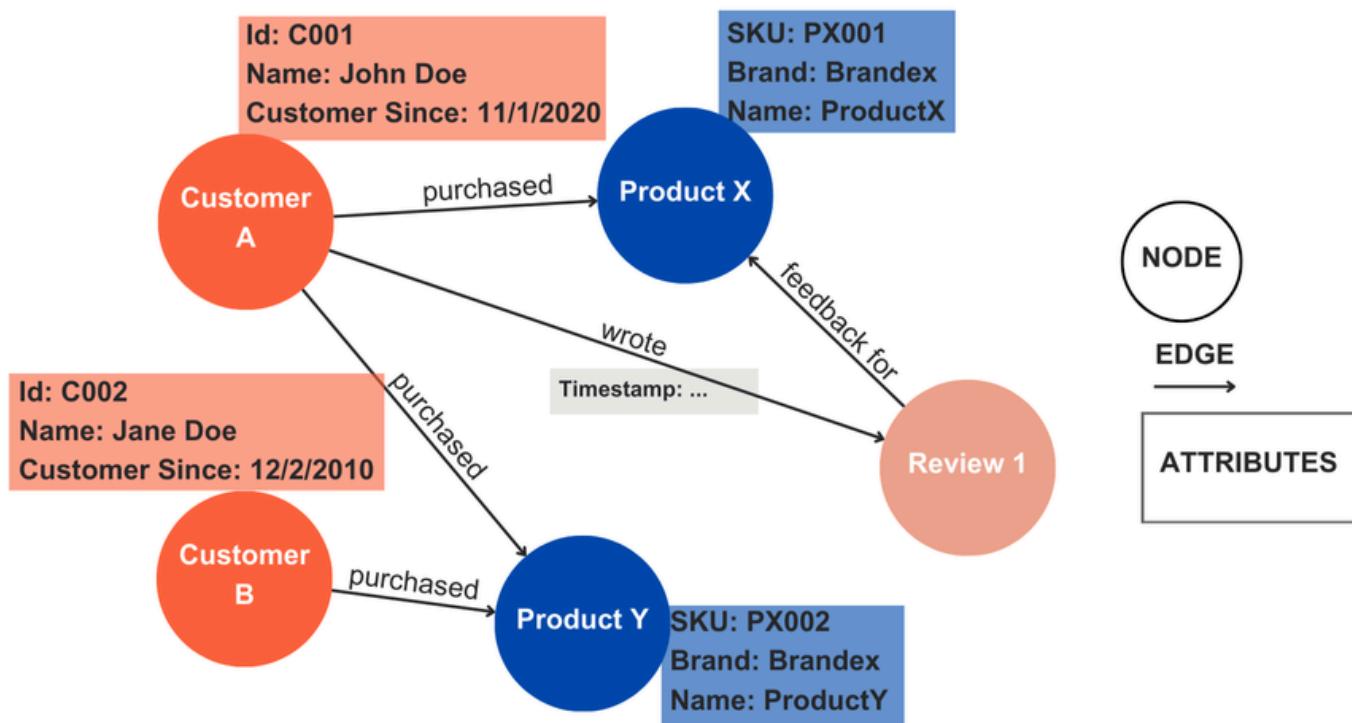
INDEXING AND GENERATION PIPELINES FOR MULTIMODAL RAG

Ch 8 : RAG variants: Multimodal, agentic, graph and other RAGs



GRAPH COMMUNITIES AND COMMUNITY SUMMARIES ENHANCED RAG

Ch 8 : RAG variants: Multimodal, agentic, graph and other RAGs



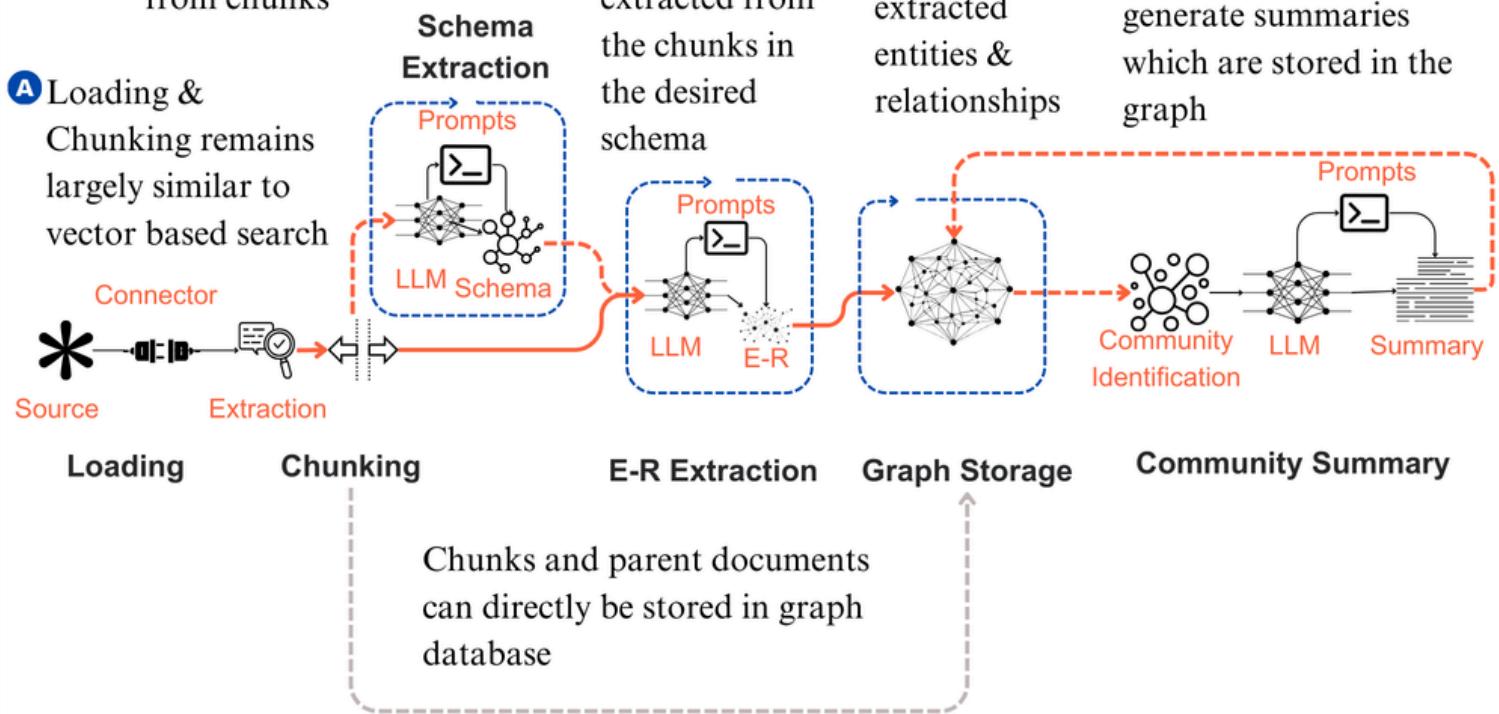
B Optionally, the graph schema can be iteratively extracted from chunks

C Entities and relationships are iteratively extracted from the chunks in the desired schema

D Graph database is updated with extracted entities & relationships

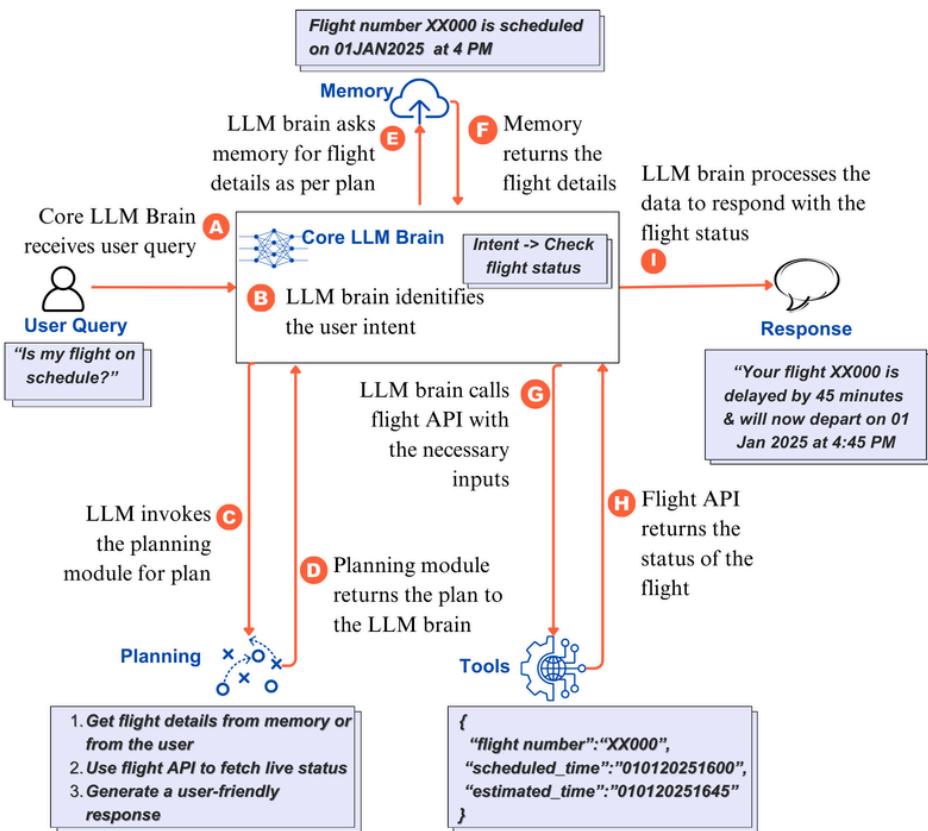
E Using a community identification algorithm, then an LLM is used to generate summaries which are stored in the graph

A Loading & Chunking remains largely similar to vector based search



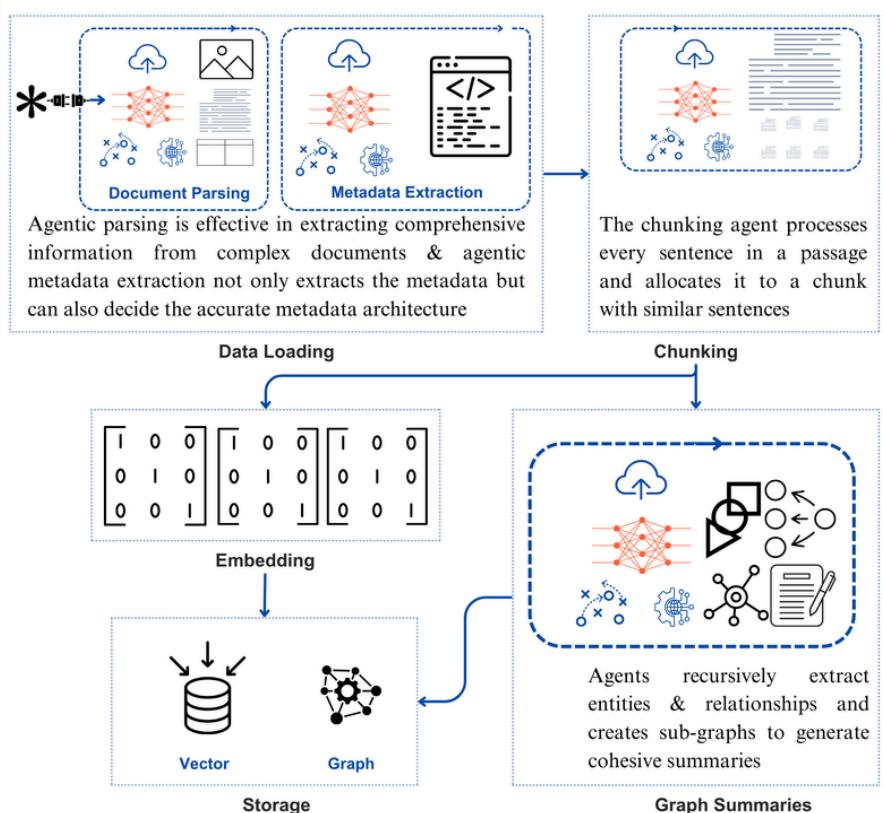
INDEXING PIPELINE FOR KNOWLEDGE GRAPH RAG

Ch 8 : RAG variants: Multimodal, agentic, graph and other RAGs



AGENTIC RAG EXAMPLE: AN AIRLINE BOOKING AND ENQUIRY SYSTEM

AGENTIC INTERVENTIONS ACROSS THE INDEXING PIPELINE STAGES



Ch 9 : RAG Development Framework & Further Exploration

This chapter covers

- A Recap of the concepts covered in this book using a six-stage RAG development framework
- Areas for further exploration in RAG



SIX STAGES OF RAG DEVELOPMENT FRAMEWORK

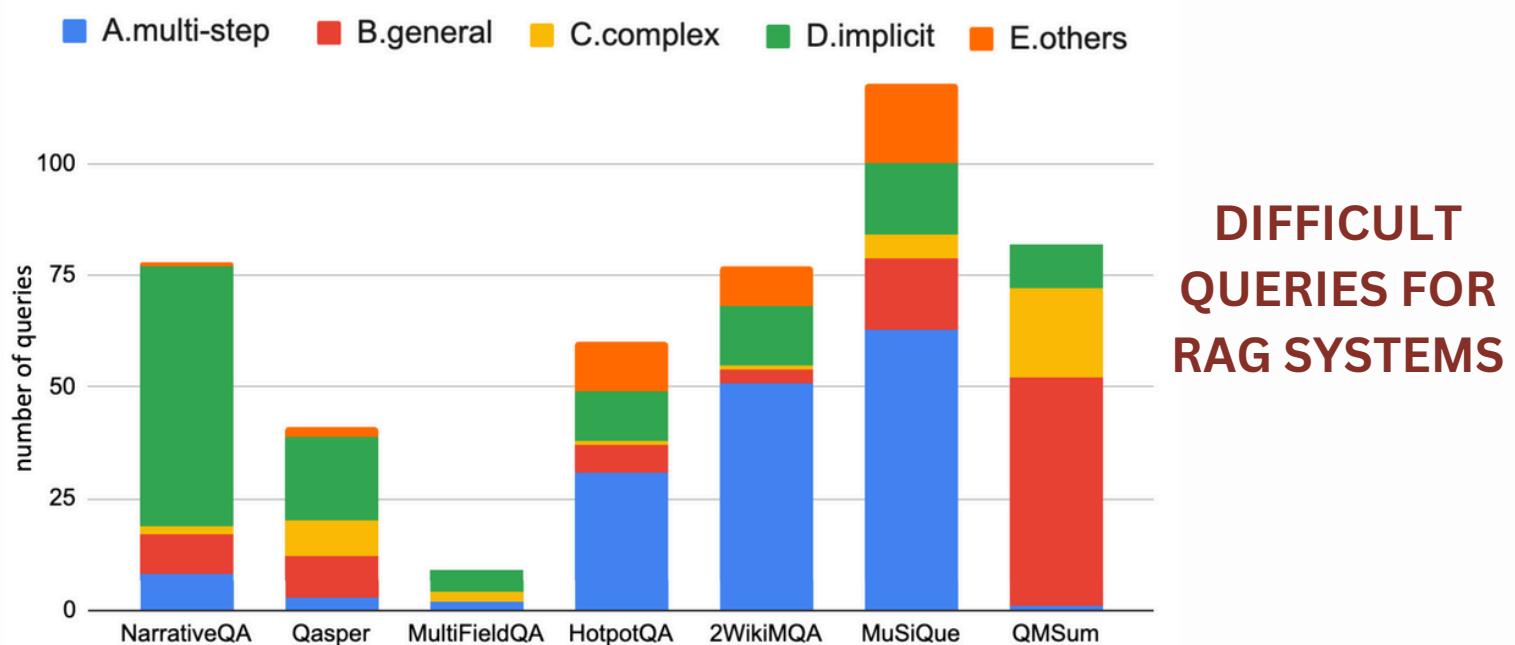
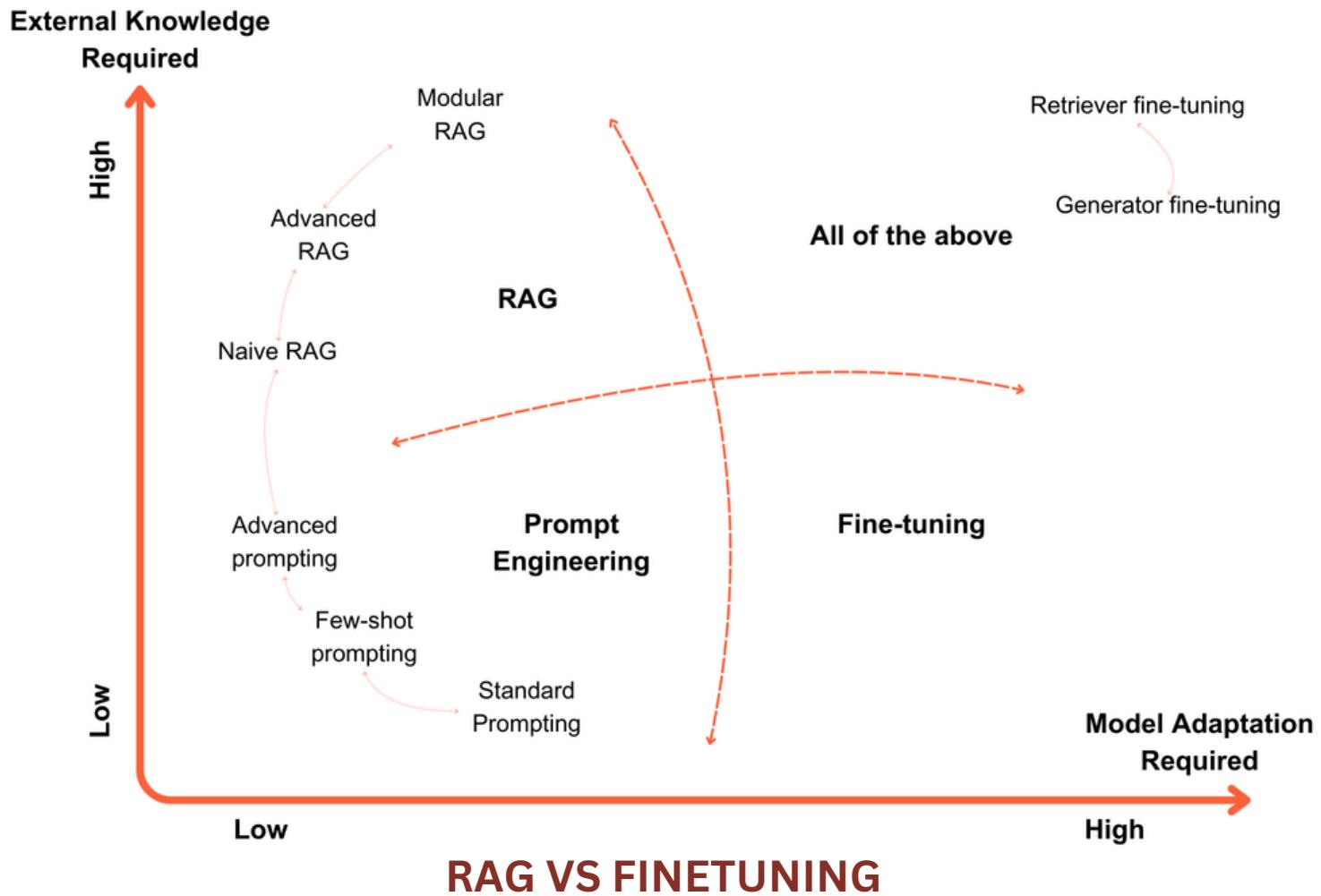
Ch 9 : RAG Development

Framework & Further Exploration

Use Case	Use Case Evaluation Questions				A RAG system is required?
	System requires data that may not be present in training set?	System requires data that is current or updates frequently?	System generates facts?	Are users looking for sources?	
Creative Writing Assistance	No - LLMs do not need any additional data for creative writing	No - LLMs do not need any current information for creative writing	Maybe - Creative writing may not necessarily need generated facts	No - User expectations from creative writing does not need any source citation	No
Customer Support Bot	Yes - Product/Company specific information may not be present in LLM training data	Yes - Product information, inventory levels, order information changes frequently	Yes - All generated information is factual	Maybe - Sources may enhance customer experience	Yes
Language Translation	No - LLMs do not need any additional data for language translation	No - LLMs do not need any current information for language translation	No - Facts, if any, will be the same as provided in the prompt	No - The source of information will always be the prompt	No
Spelling & Grammar Correction	No - LLMs do not need any additional data for checks	No - LLMs do not need any current information for checks	No - No additional facts need to be generated	No - The source of information will always be the prompt	No

EVALUATING USE CASES FOR THE NEED OF A RAG SYSTEM

Ch 9 : RAG Development Framework & Further Exploration

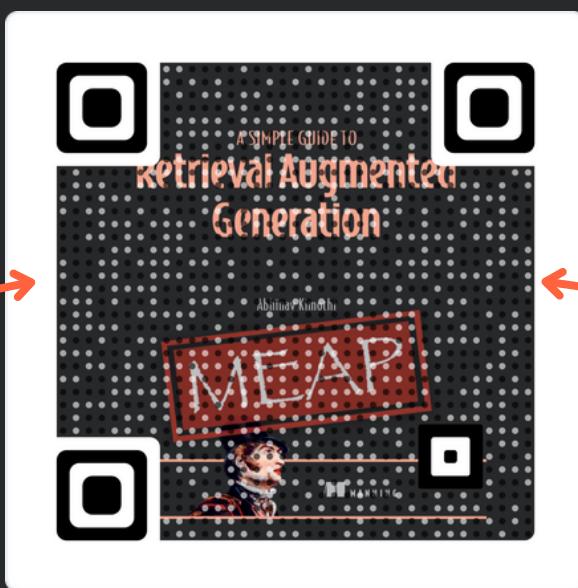
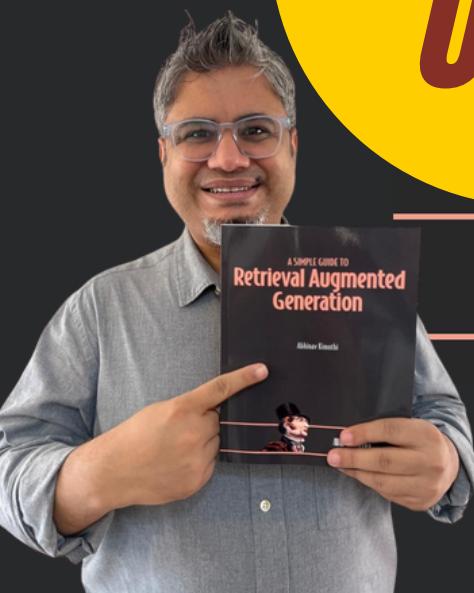


THESE EXCERPTS ARE FROM A SIMPLE GUIDE TO RETRIEVAL AUGMENTED GENERATION

A SIMPLE GUIDE TO
**Retrieval Augmented
Generation**

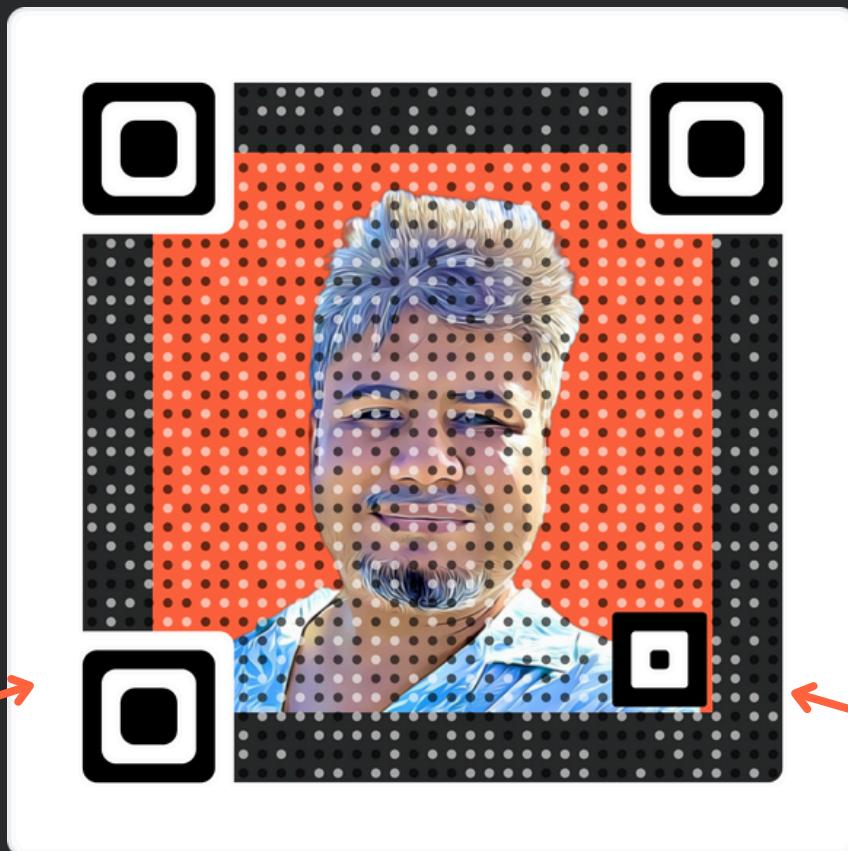
**50%
OFF**

1400+
Copies Sold



**SCAN CODE OR CLICK HERE
TO GET AN EARLY ACCESS COPY**

Hi! My name is Abhinav, and I talk about ML, RAG, LLMs & AI Agents. If our interests align, I'd love to stay connected



**SCAN CODE OR CLICK HERE
TO CONNECT**



linktr.ee/abhinavkimothi



[/in/Abhinav-Kimothi](https://www.linkedin.com/in/Abhinav-Kimothi)



[@akaiworks](https://www.instagram.com/@akaiworks)



[@abhinav_kimothi](https://twitter.com/abhinav_kimothi)



[@abhinavkimothi](https://discord.com/users/100000000000000000)