

# Retrieval Augmented Generation



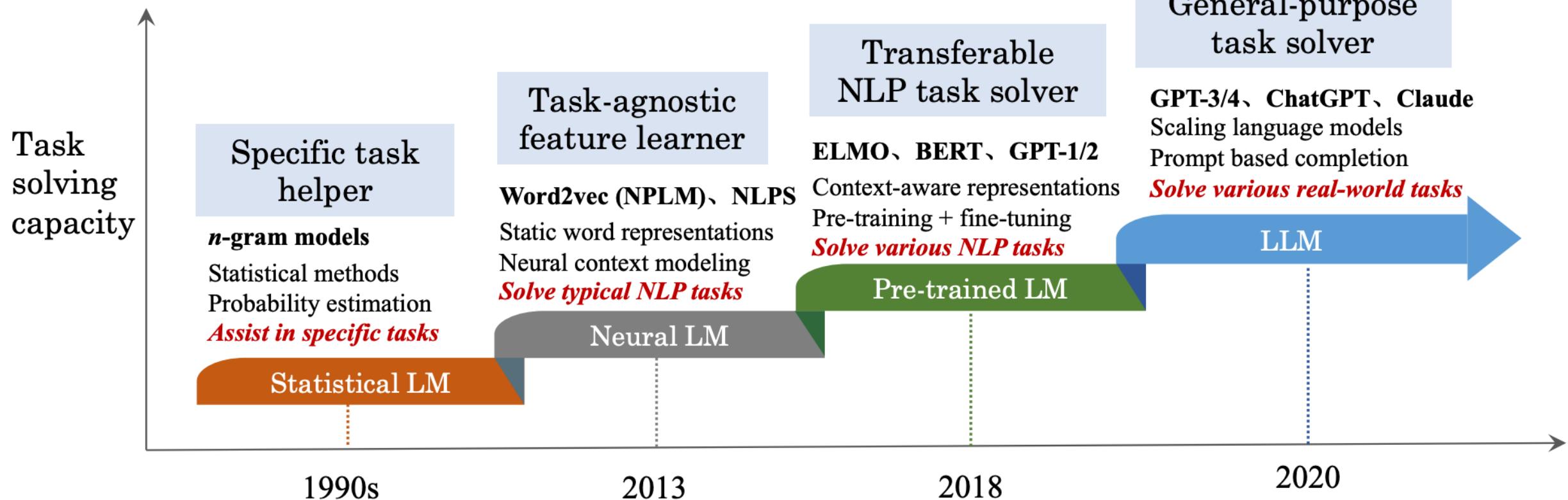
INSTITUT  
POLYTECHNIQUE  
DE PARIS

Mehwish Alam  
Associate Professor  
Télécom Paris  
Institut Polytechnique de Paris  
Winter Semester 2024-2025

The background of the image is a complex network graph. It consists of numerous small, semi-transparent blue and grey circular nodes of varying sizes scattered across the frame. These nodes are interconnected by a dense web of thin, light blue lines representing edges. Some larger, more prominent white circles are also present, particularly towards the center and bottom left, which appear to be hubs or focal points for clusters of smaller nodes.

# Brief Recap

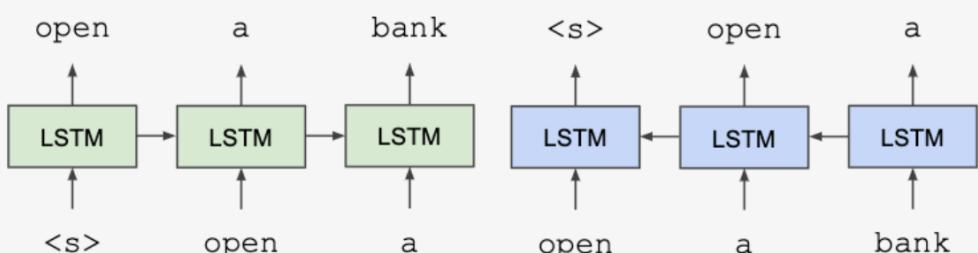
# Evolution of Language Models



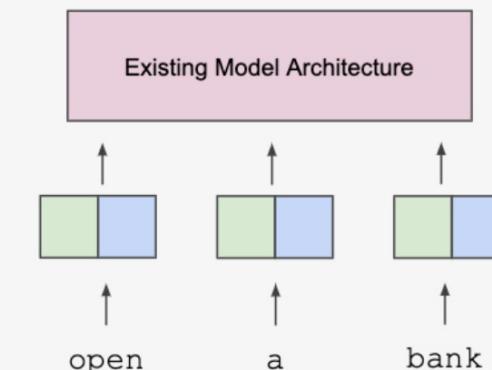
# Prior Work: ELMo

- ELMo (Peters et al., 2018; NAACL 2018 best paper)
  - Train **two separate unidirectional LMs** (left-to-right and right-to-left) based on LSTMs
  - Feature-based approach: **pre-trained representations** used as input to task-specific models
  - Trained on single sentences from 1B word benchmark ([Chelba et al., 2014](#))

## Train Separate Left-to-Right and Right-to-Left LMs



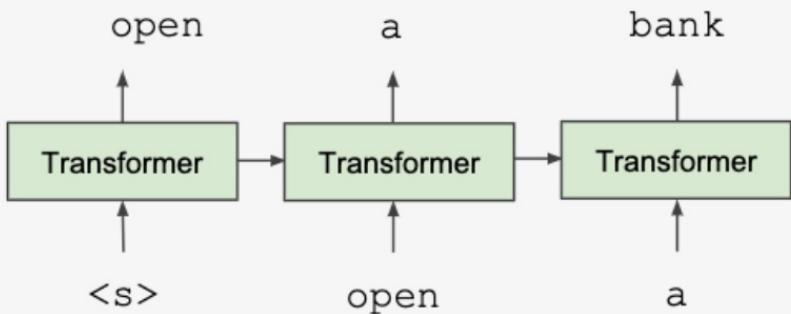
## Apply as “Pre-trained Embeddings”



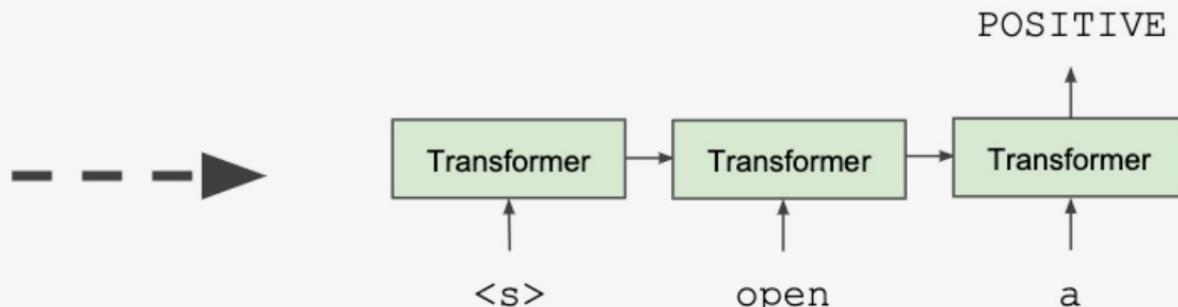
# Prior Work: OpenAI GPT

- OpenAI GPT (Radford et al., 2018; released in 2018/6)
  - Train **one unidirectional LM (left-to-right)** based on a deep Transformer decoder
  - Fine-tuning approach: all pre-trained parameters are re-used & updated on downstream tasks
  - Trained on **512-token segments on BooksCorpus** — much longer context!

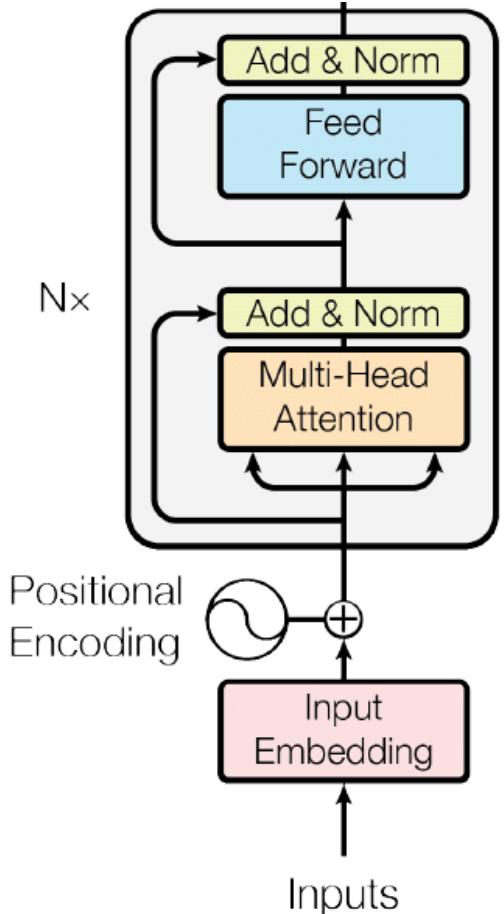
## Train Deep (12-layer) Transformer LM



## Fine-tune on Classification Task



# BERT pre-training: putting together

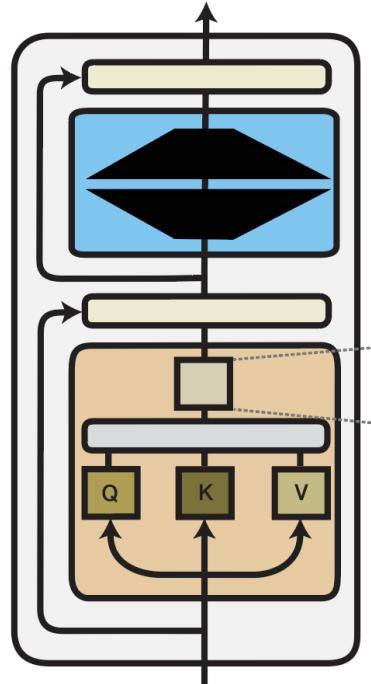


- BERT-base: 12 layers, 768 hidden size, 12 attention heads, 110M parameters
- BERT-large: 24 layers, 1024 hidden size, 16 attention heads, 340M parameters
- Training corpus: Wikipedia (2.5B) + BooksCorpus (0.8B)
- Max sequence size: 512 word pieces (roughly 256 and 256 for two non-contiguous sequences)
- Trained for 1M steps, batch size 128k

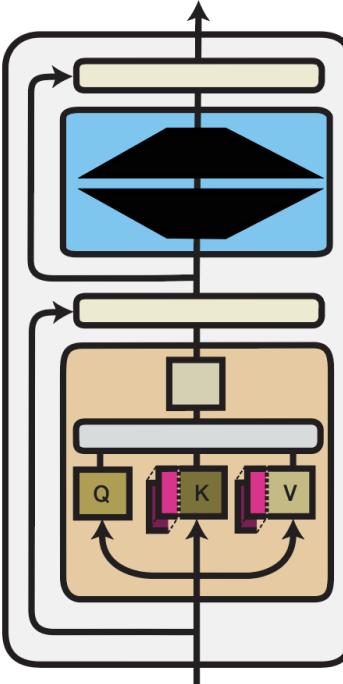
Same as OpenAI GPT

OpenAI GPT was trained on BooksCorpus only!

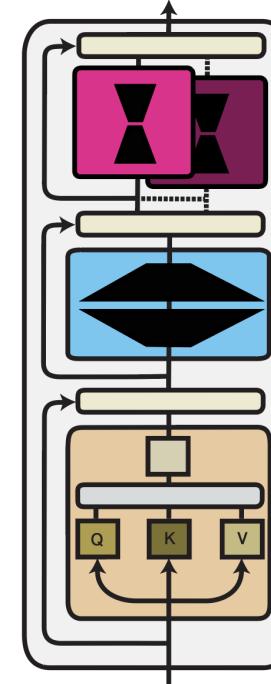
# Three Computation Functions



Parameter Composition



Input Composition

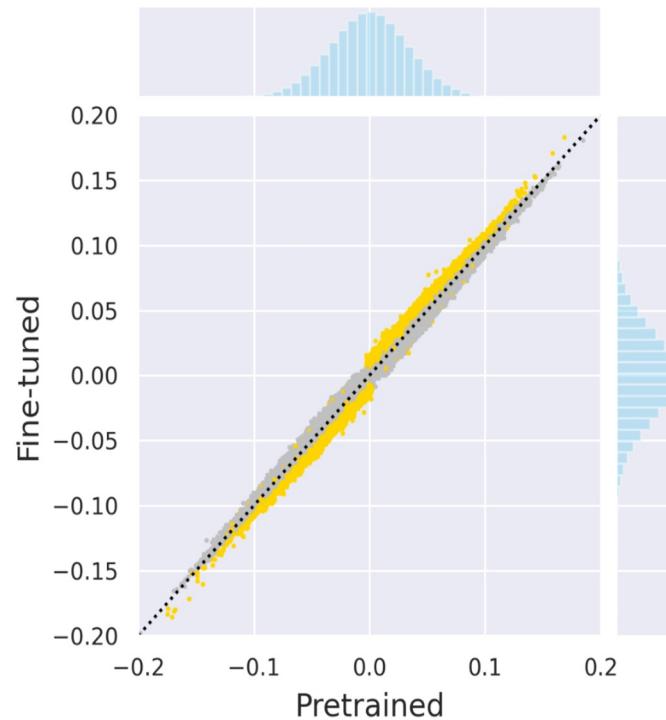
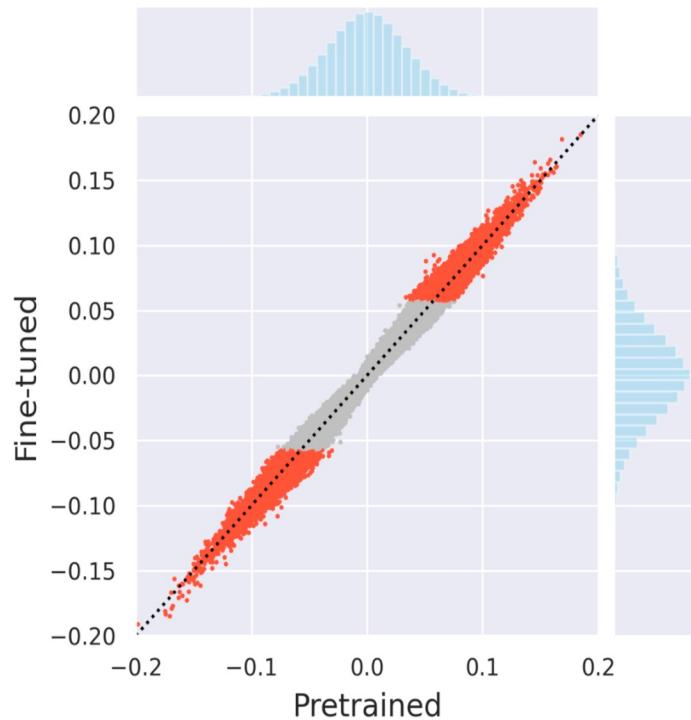


Function Composition

# Pruning Pre-trained Models

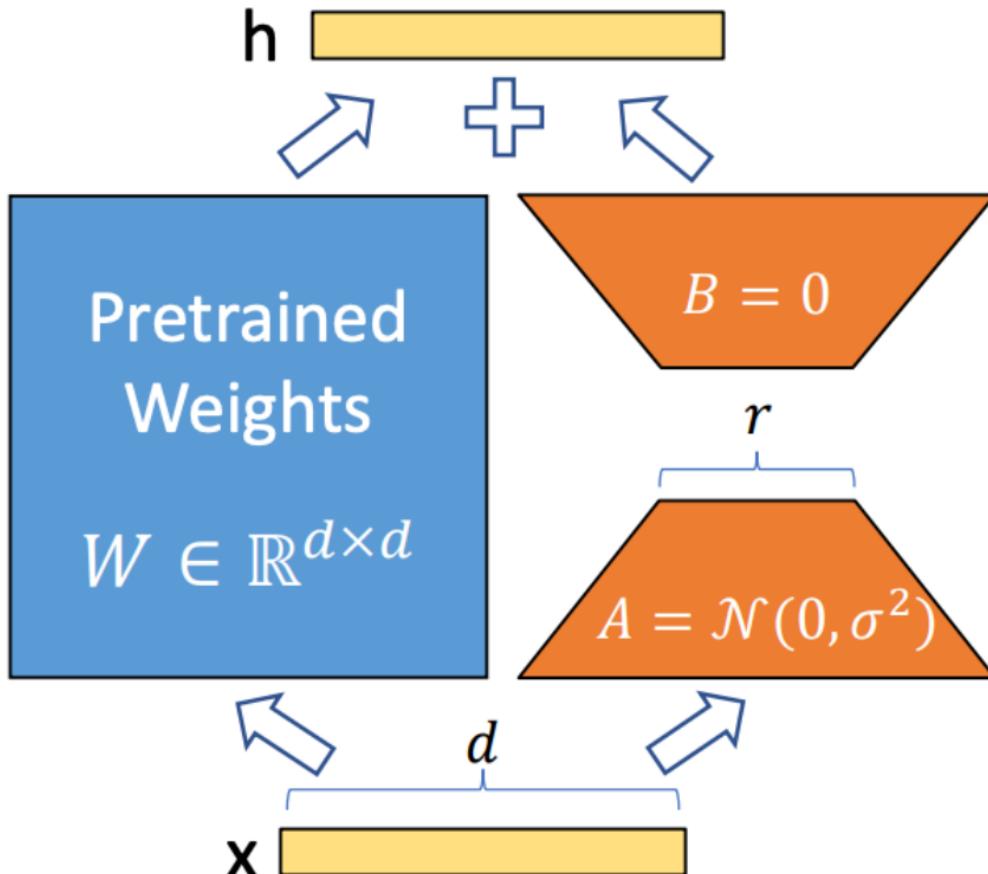
- Pruning does not consider how weights change during fine-tuning
- **Magnitude pruning:** keep weights farthest from 0
- **Movement pruning [Sanh et al., 2020]:** keep weights that *move the most away* from 0

Fine-tuned weights stay close to their pre-trained values.  
Magnitude pruning (left) selects **weights that are far from 0.**

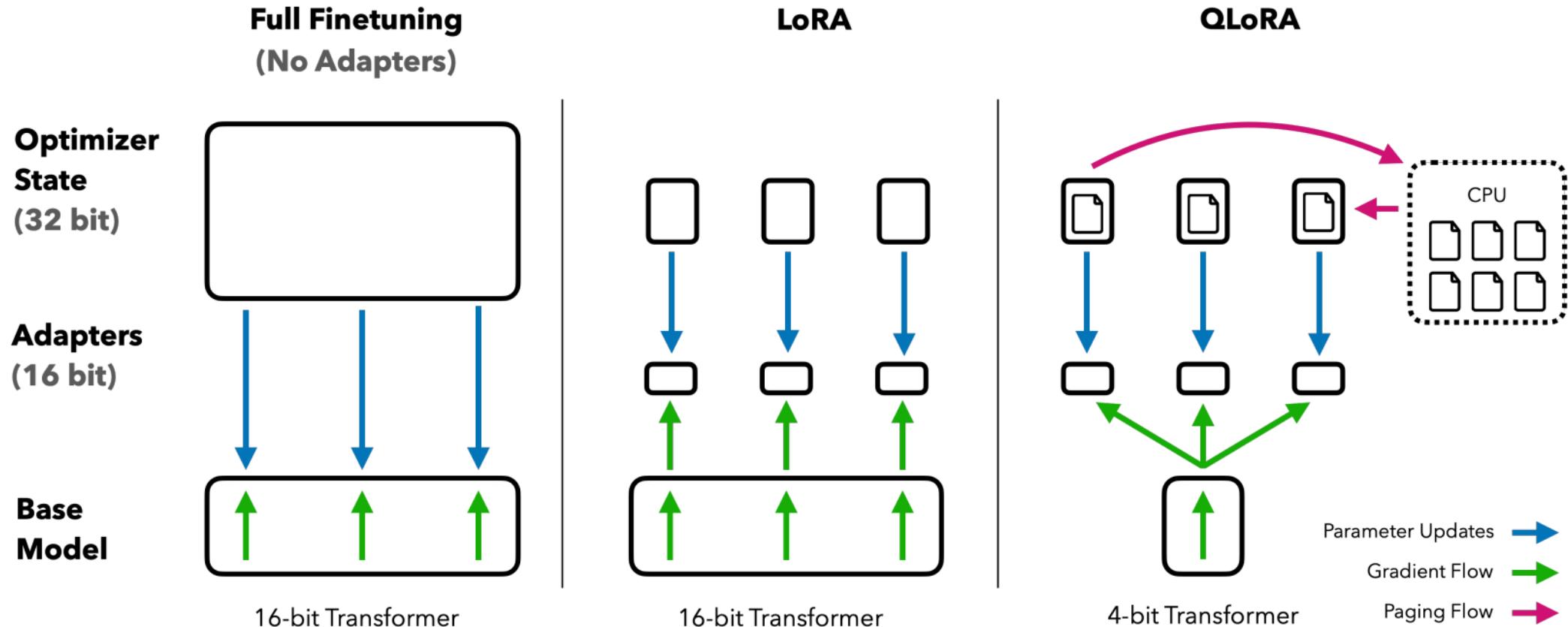


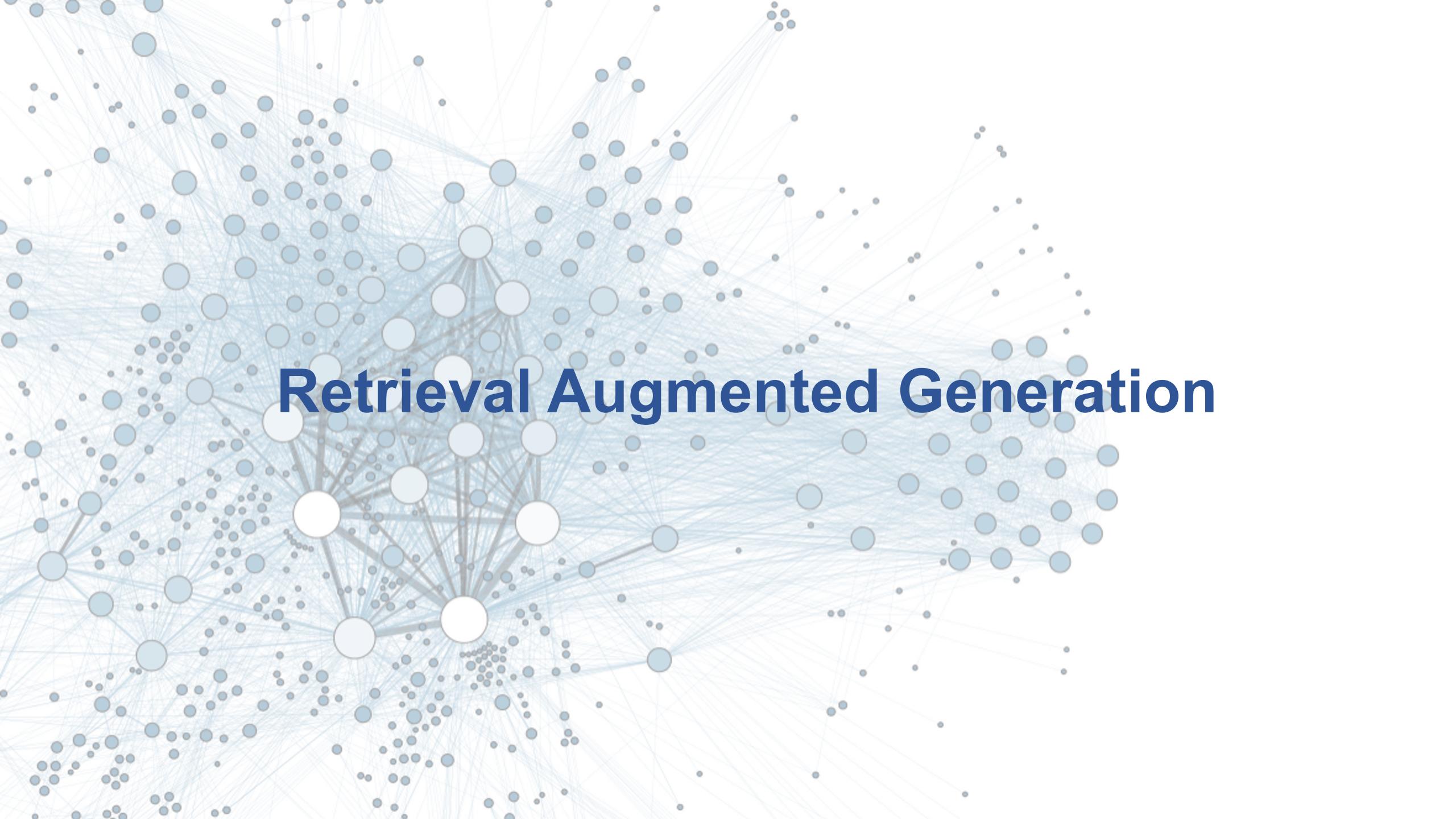
Movement pruning (right) selects weights that **move away from 0.**

# Low-rank-parameterized update matrices



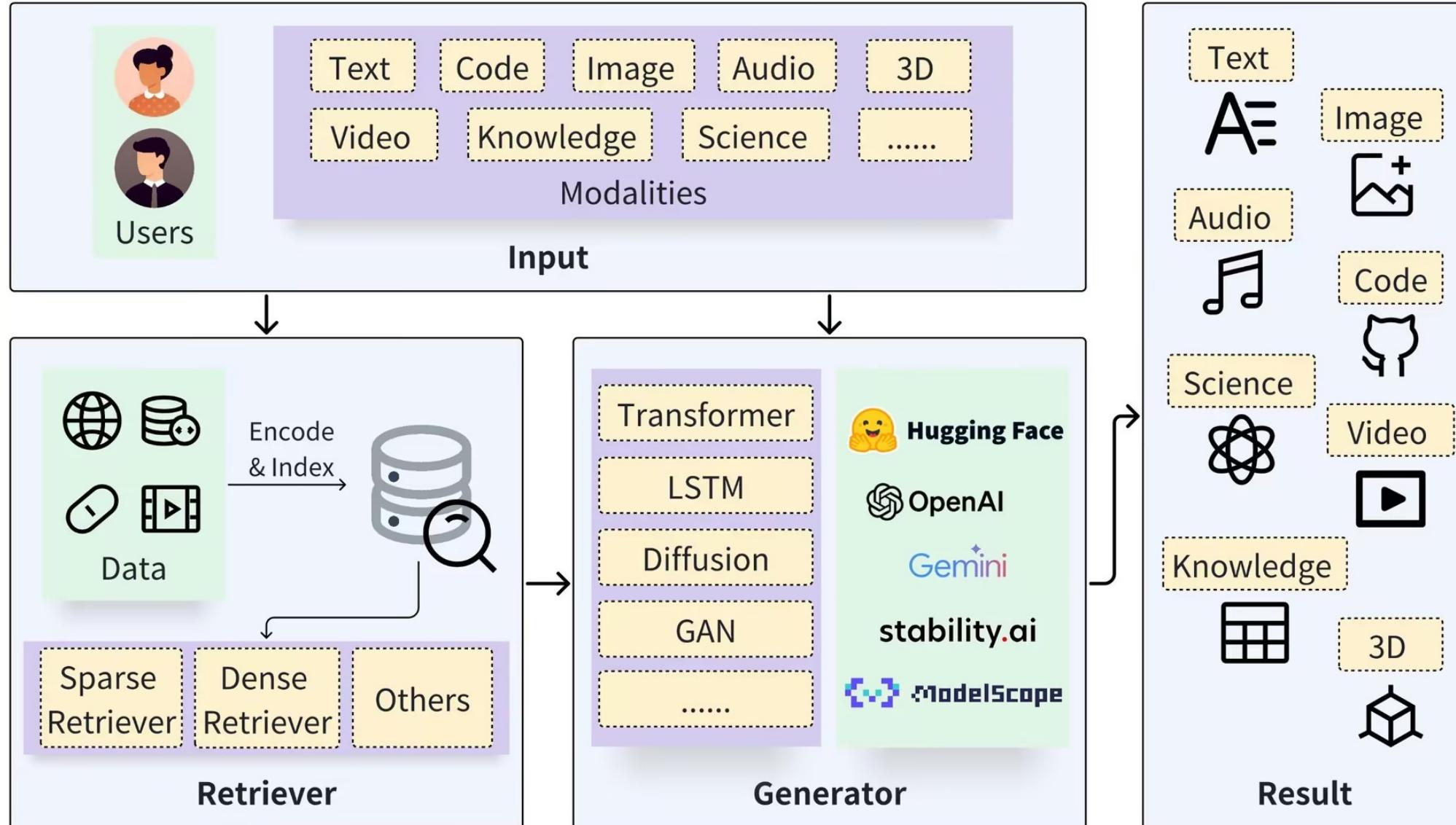
# From LoRA to QLoRA



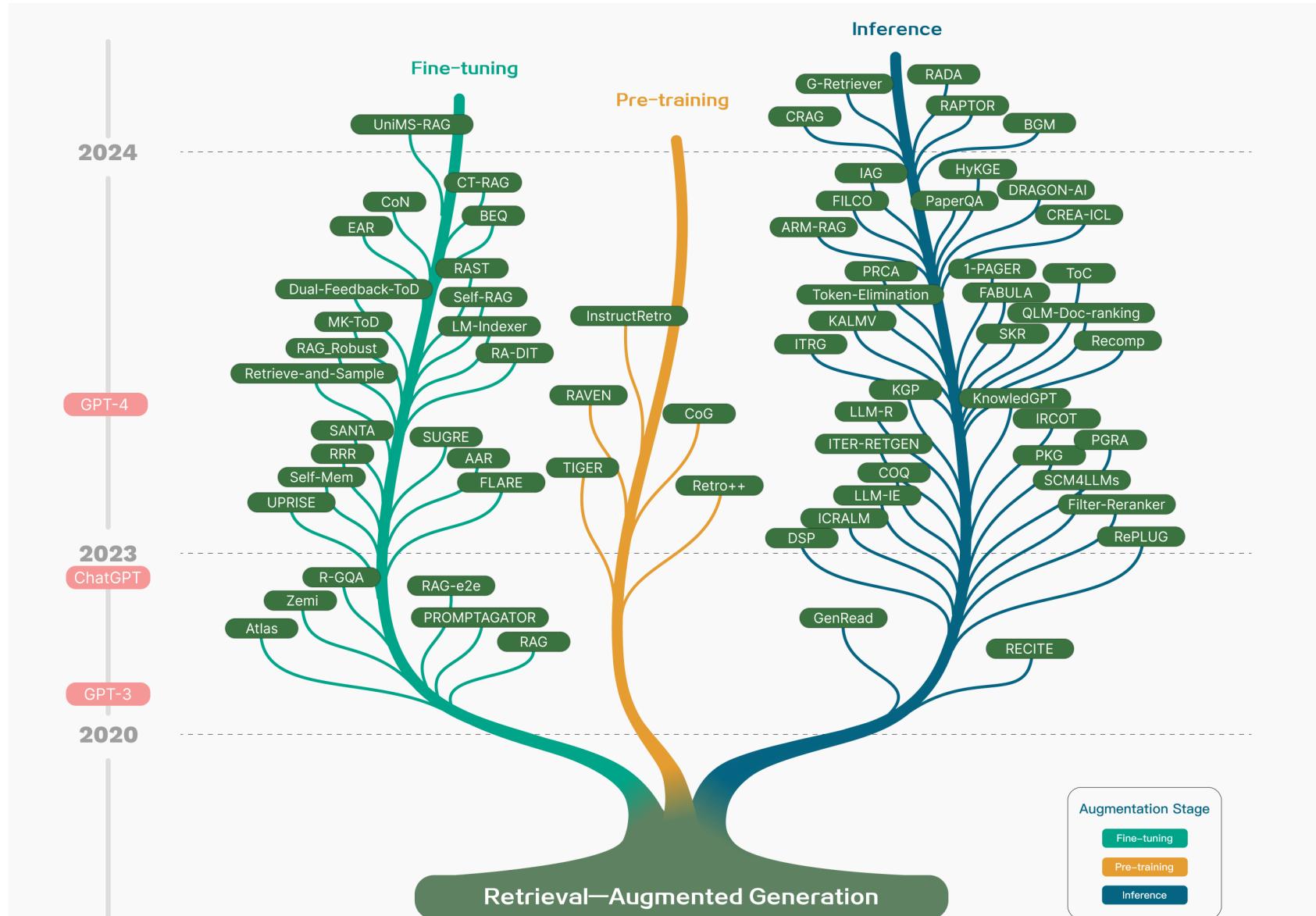


# Retrieval Augmented Generation

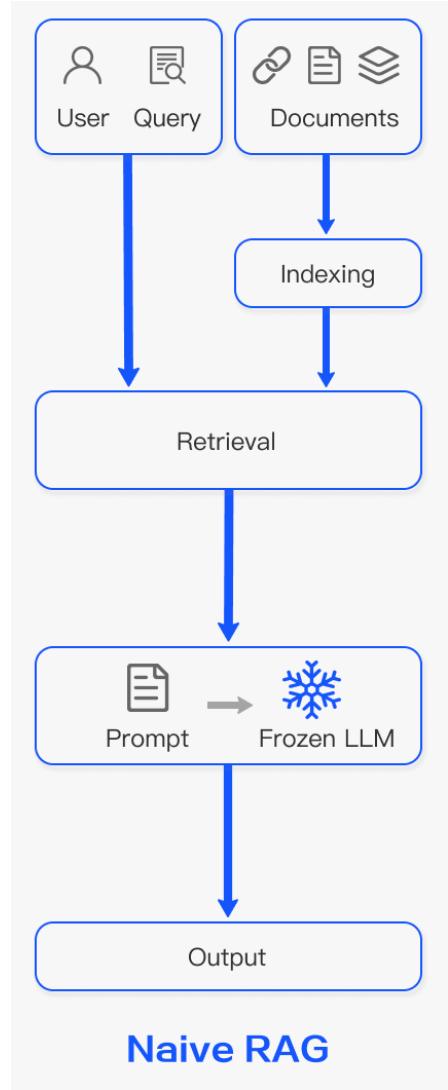
# Retrieval Augmented Generation



# Retrieval-Augmented Generation

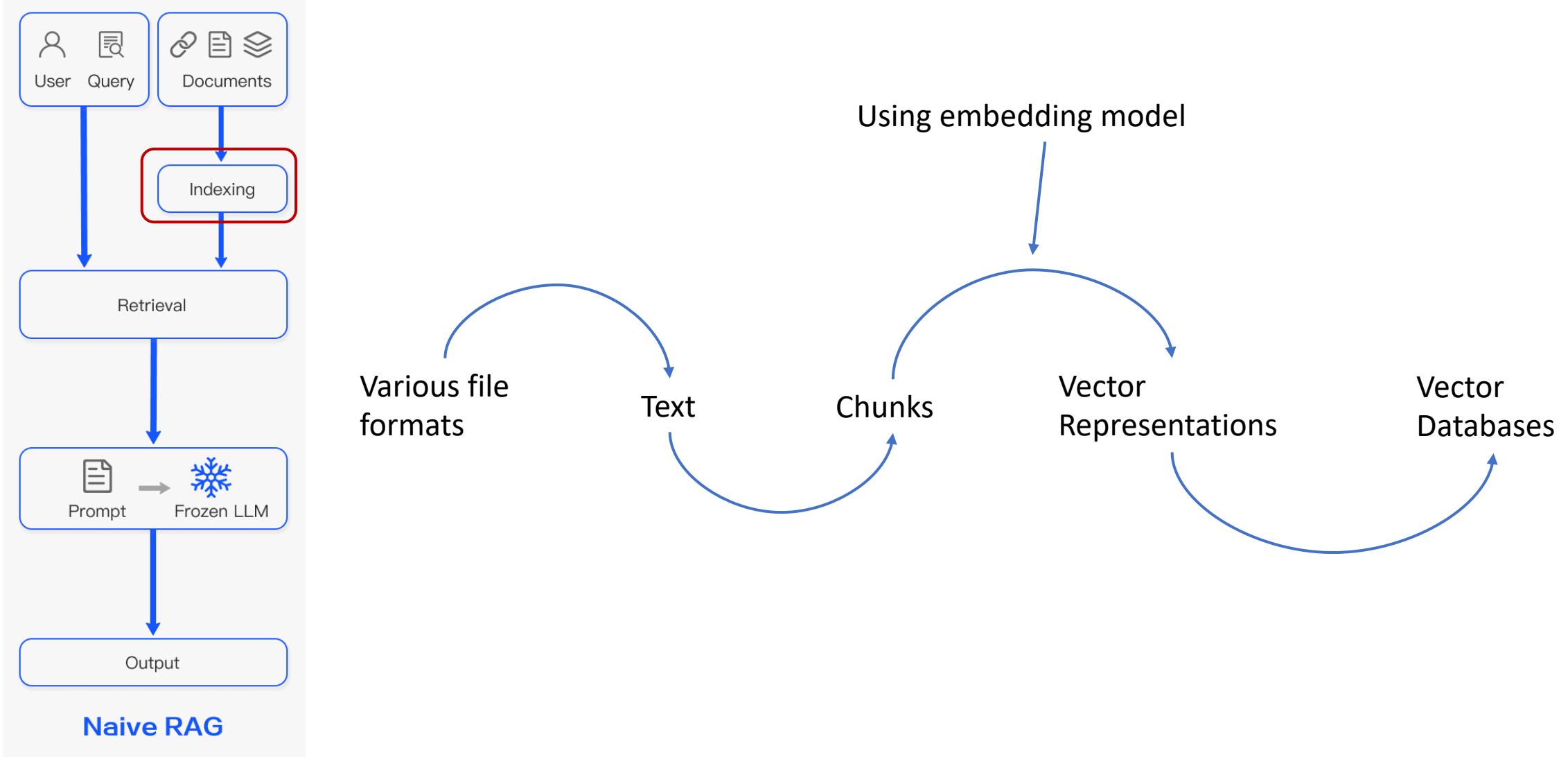


# Naive RAG

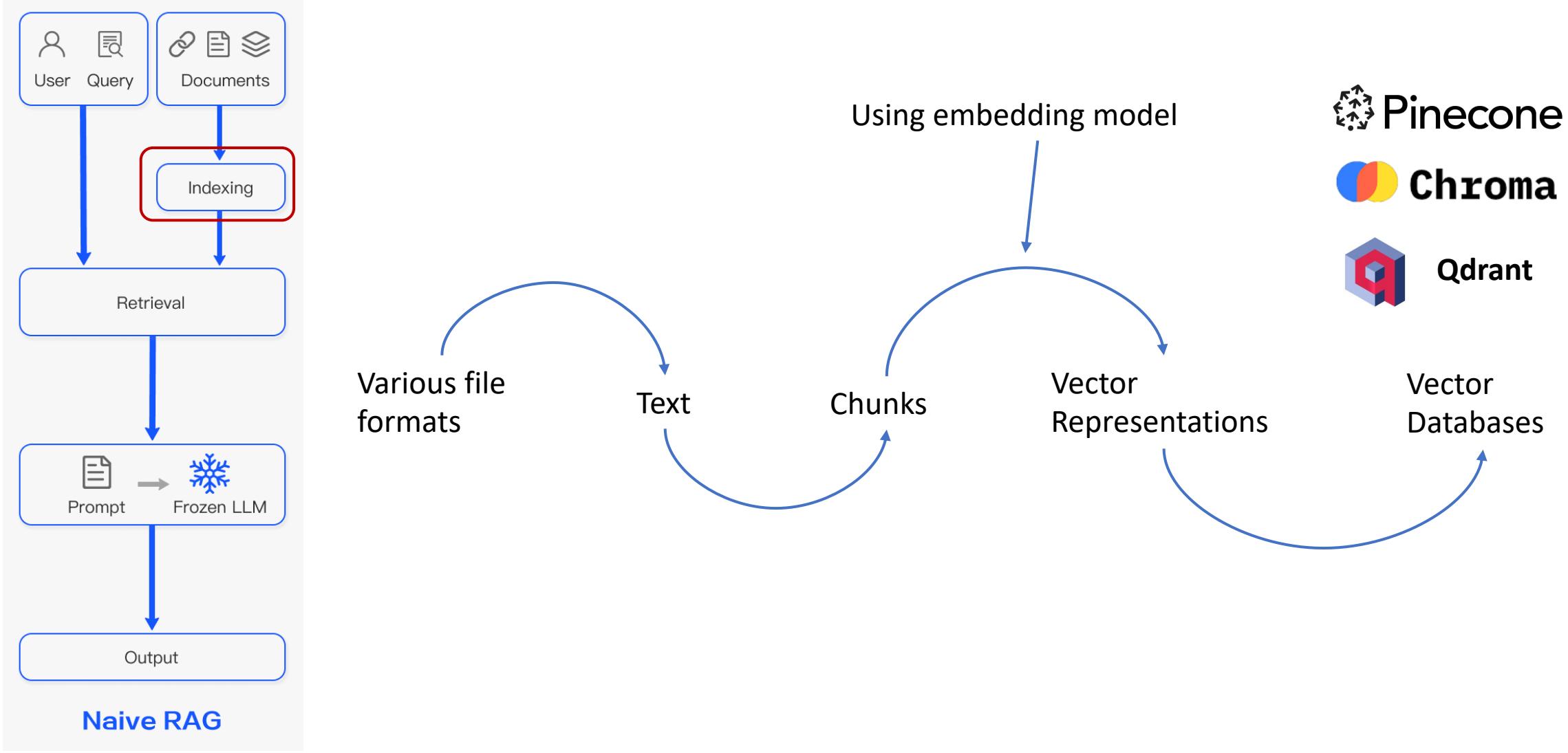


Retrieve Read Framework

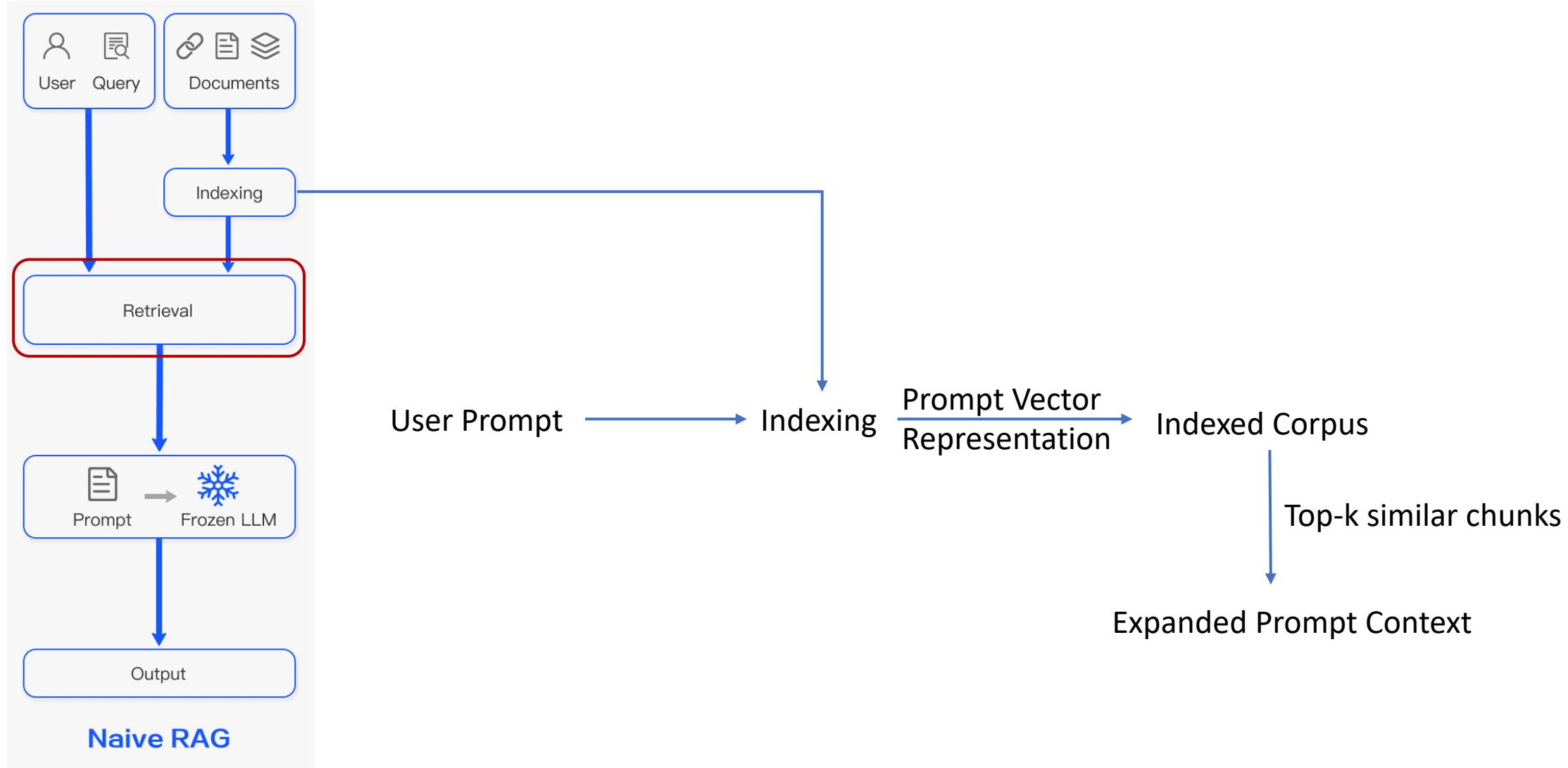
# Naive RAG - Indexing



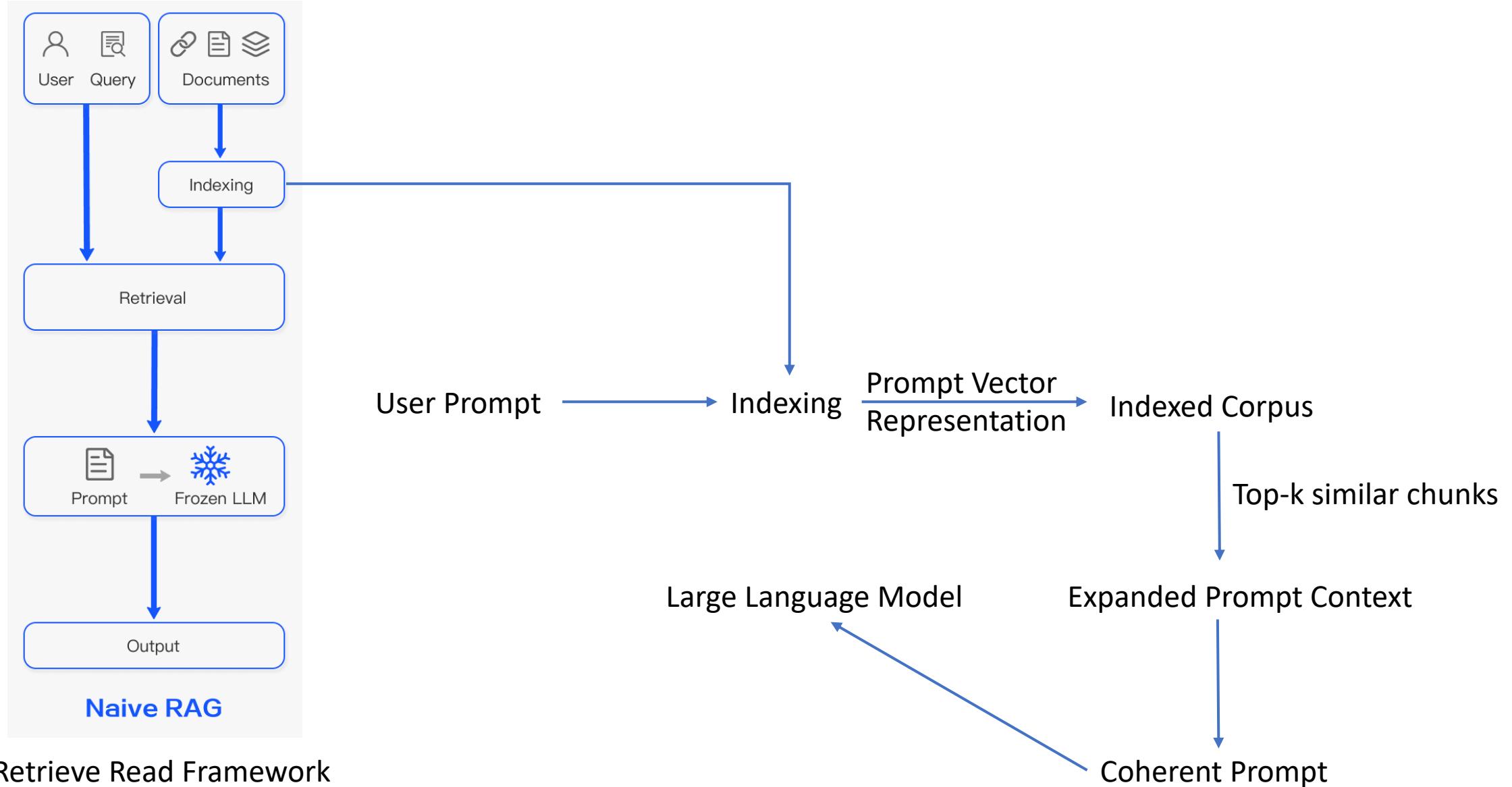
# Naive RAG - Indexing



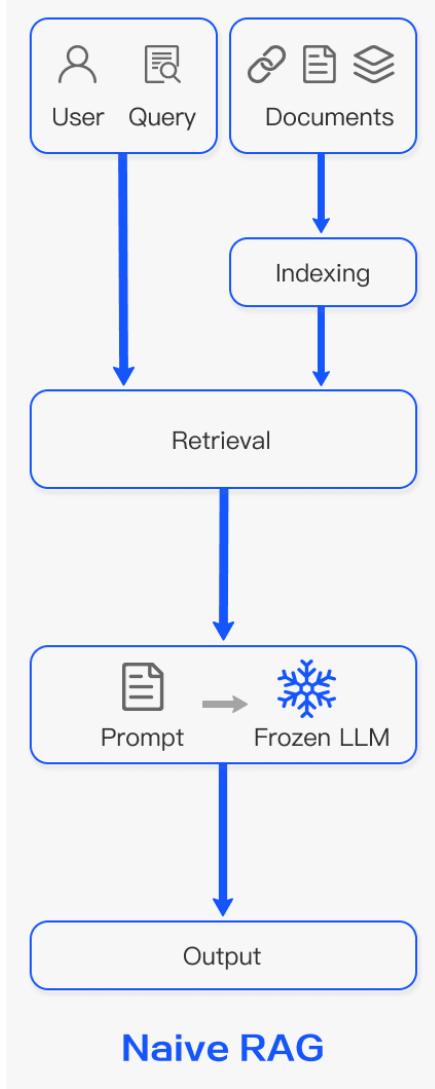
# Naive RAG - Retrieval



# Naive RAG - Generation



# Naive RAG - Drawbacks



## Retrieval.

- Struggles with precision and recall:
  - Selection of **misaligned or irrelevant** chunks
  - **Missing** of crucial information.

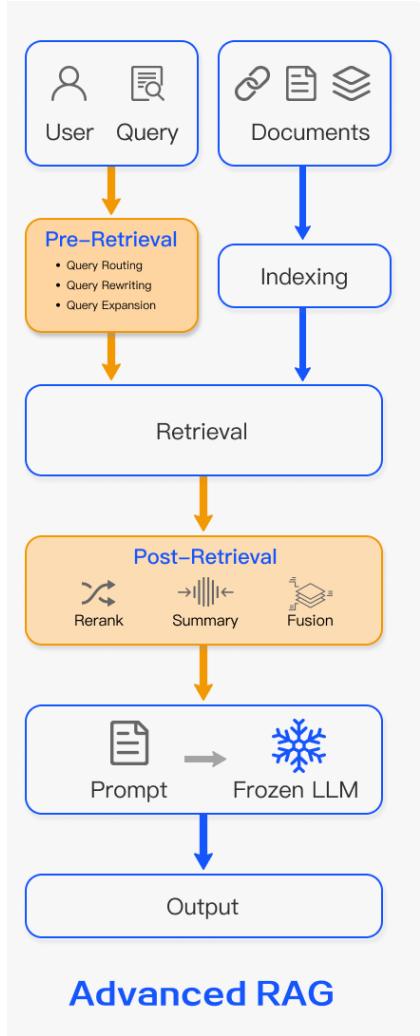
## Generation.

- Hallucinations
- Irrelevance
- Toxicity
- Bias

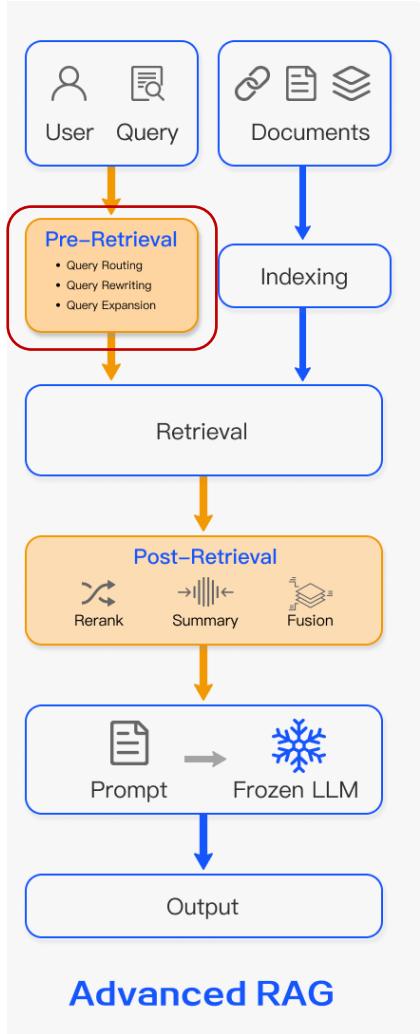
## Augmentation.

- Integrating retrieved information
- Redundancy
- Over reliance on augmented information

# Advanced RAG



# Advanced RAG



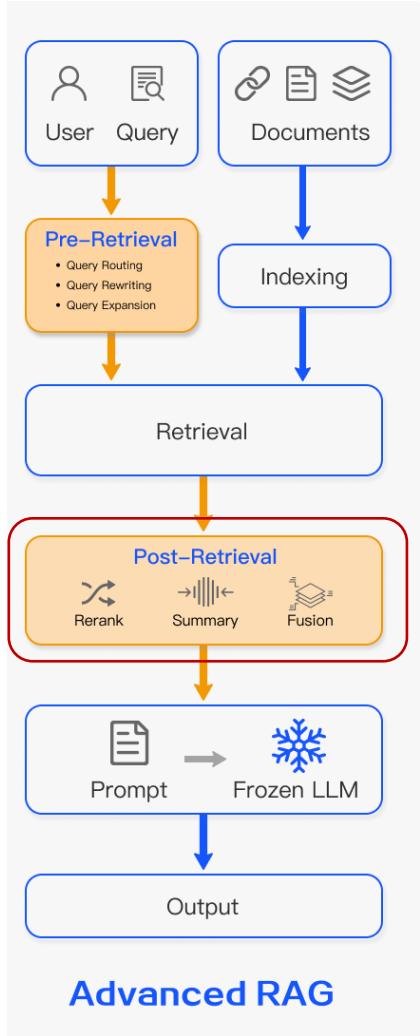
## Pre-retrieval Process.

- Optimizing the indexing structure and the original query.
- Enhance the quality of the content being indexed
  - Enhancing data granularity
  - optimizing index structures
  - adding metadata
  - alignment optimization
- Query optimization

Original question is clearer and more suitable for the retrieval task.

- Query transformation
- Query rewriting
- Query expansion, etc.

# Advanced RAG



## Post-Retrieval Process.

Integrate retrieved context effectively with the query.

- Re-ranking
- Implementations: LlamaIndex, LangChain, and HayStack

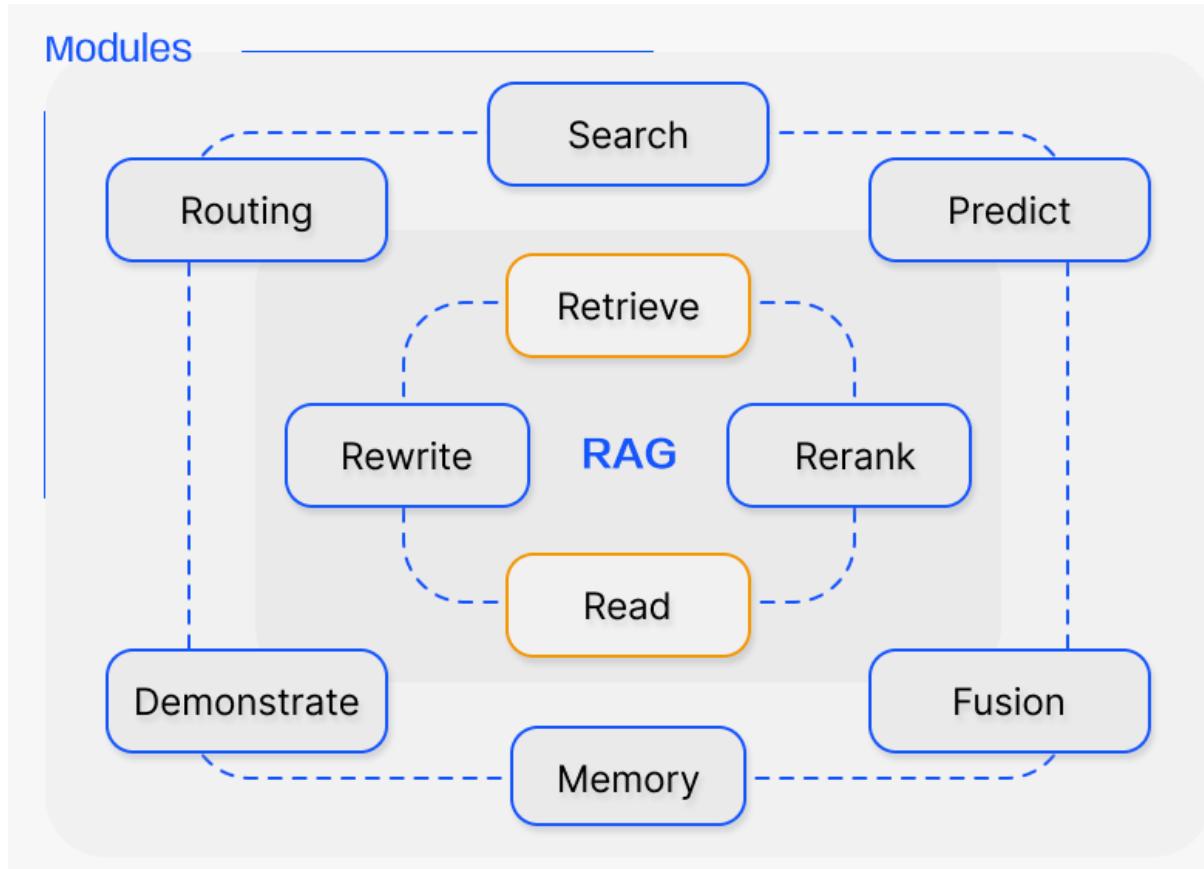
Drawback.

- Information overload
- Irrelevant content

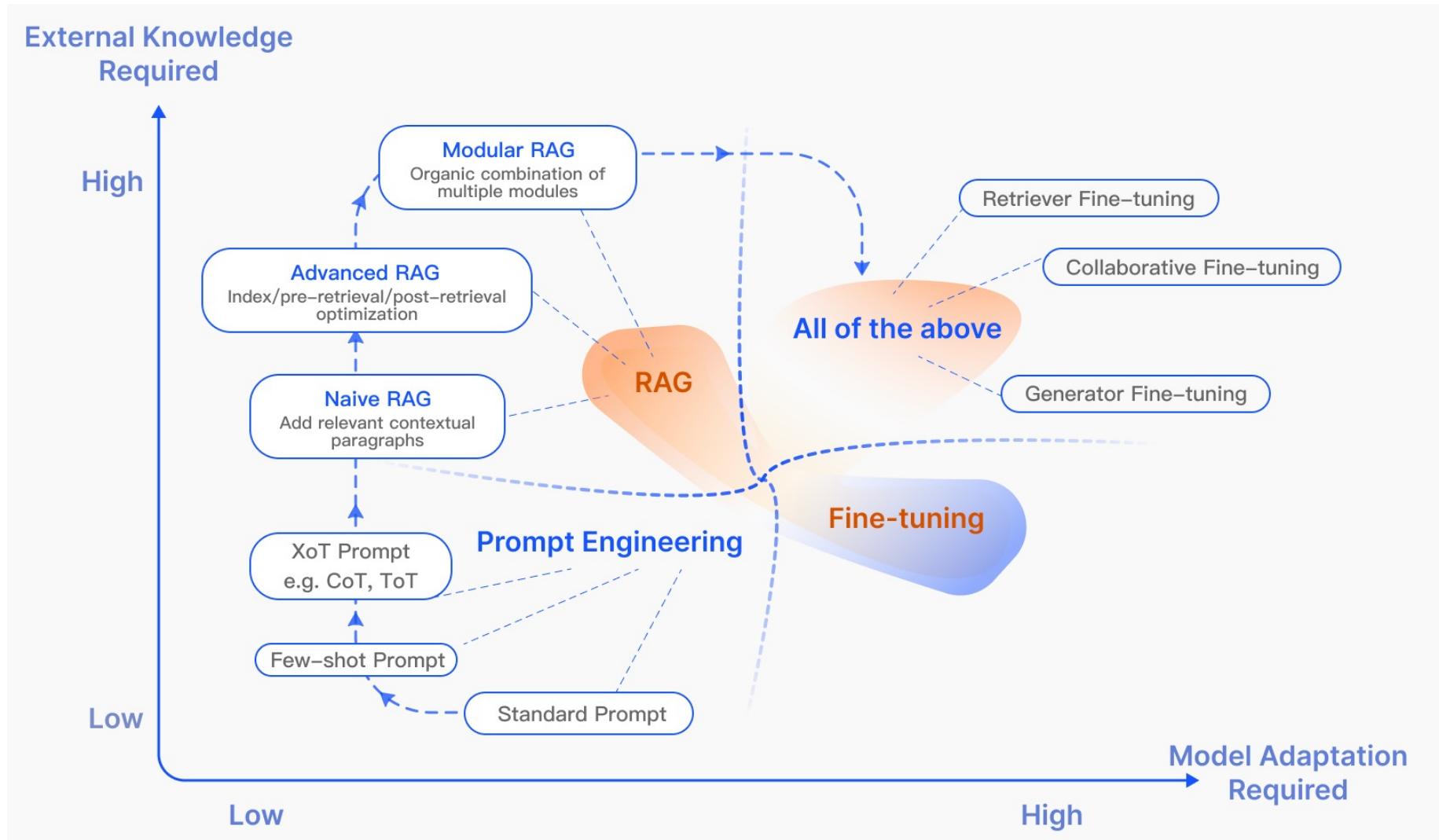
Mitigation.

- Selecting the essential information, emphasizing critical sections
- Shortening the context to be processed

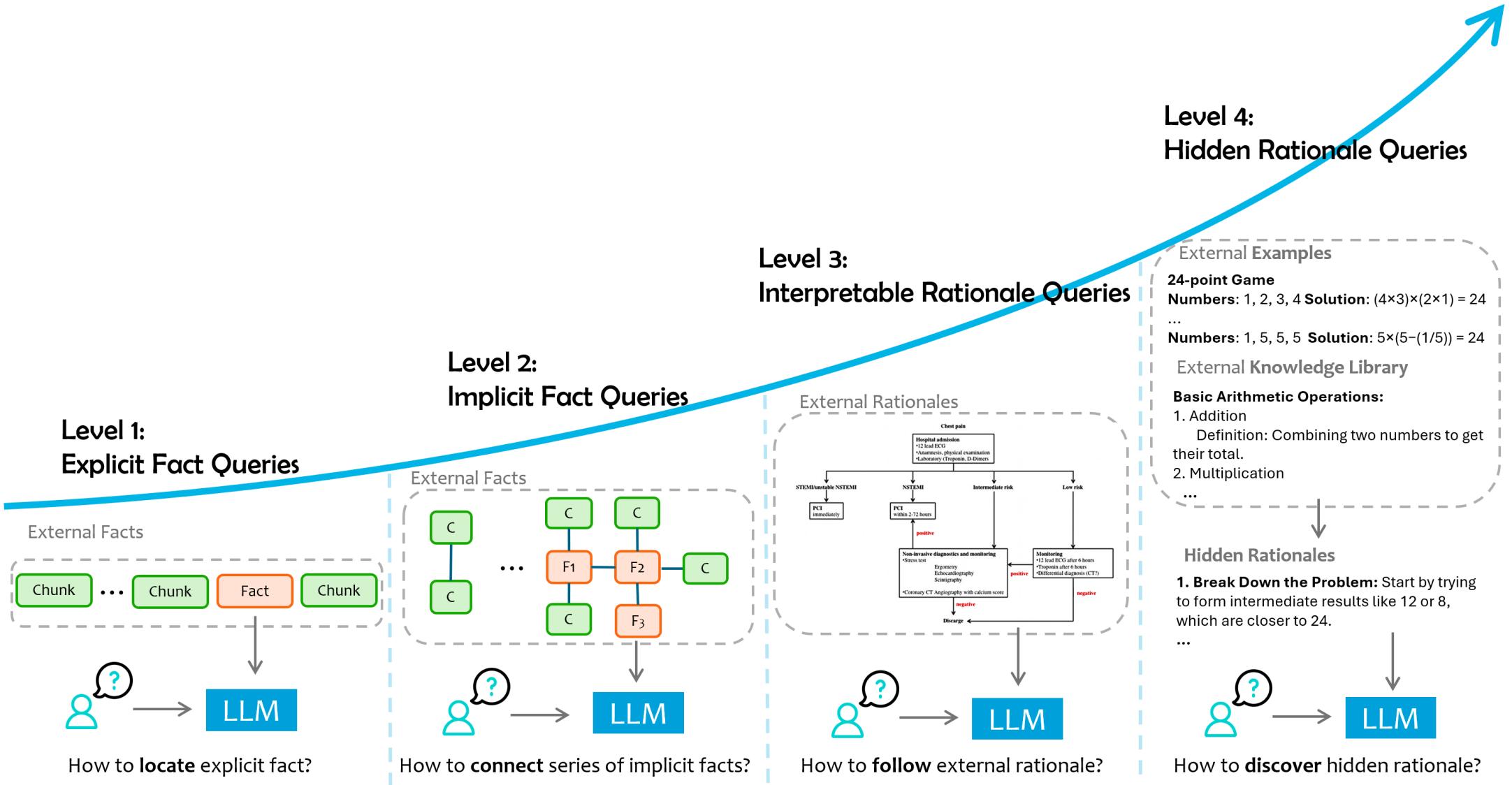
# Modular RAG



# RAG - Comparison



# Levels of Queries



# Level 1 – Explicit Facts

---

- Explicit facts directly present in the domain specific data
- Do not require any additional reasoning

## Example.

“Where were the Summer Olympics 2024 held?”

“When did Keith Flint passed away?”

# Level 1 – Explicit Facts

---

## Challaneges and Solutions

### Data Processing Difficulties

- Data is often highly unstructured
- contains **multi-modal components** such as tables, images, videos, and more
- Segmentation or "chunking" poses **difficulties in retaining the original context and meaning**

### Data Retrieval Difficulties

- Retrieval of relevant data segments from a large, unstructured dataset can be **computationally intensive** and **prone to errors**

### Evaluation Difficulties

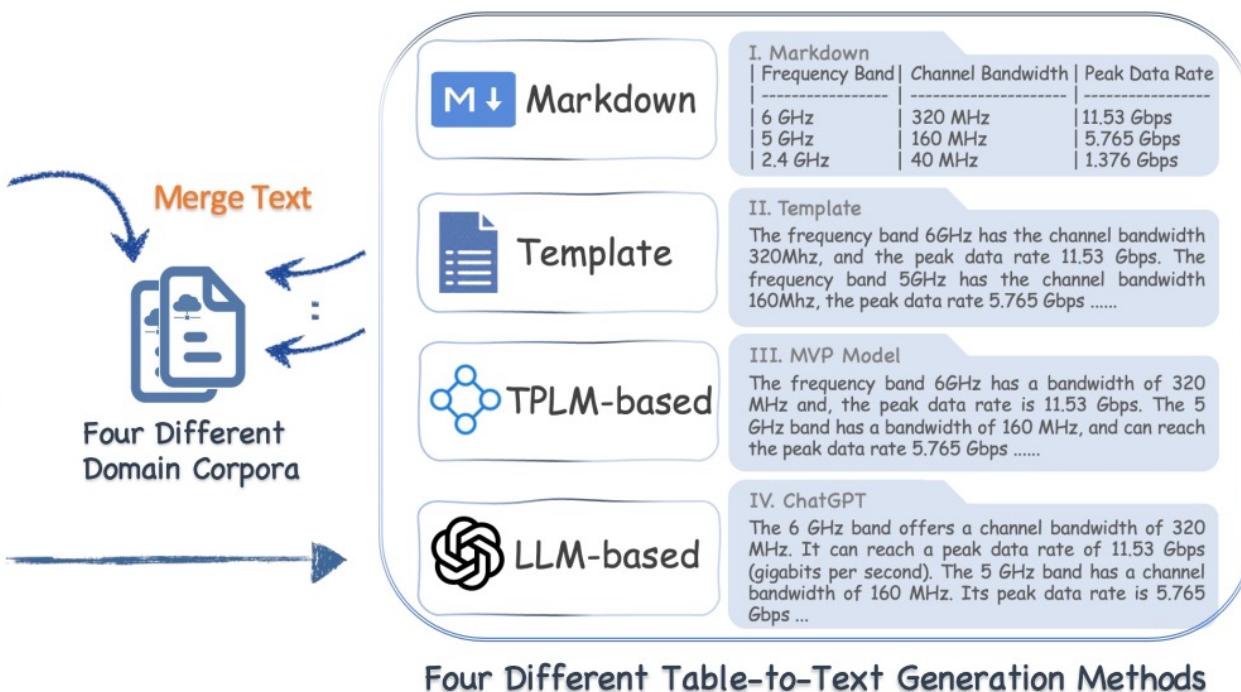
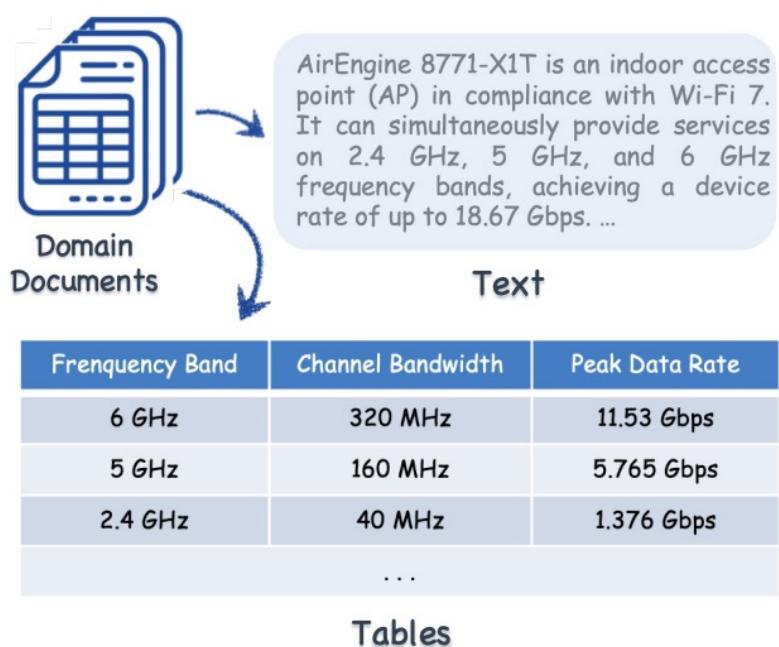
- Requires robust **metrics** for assessing the quality of data retrieval and response generation

# Level 1 – Explicit Facts - RAG

## Data Processing Enhancement

### Multimodal Documents Parsing

- Converting multi-modal content into textual form
- Multimodal embeddings as soft prompts for input



# Level 1 – Explicit Facts - RAG

## Data Processing Enhancement

### Multimodal Documents Parsing

- Converting multi-modal content into textual form
- Multimodal embeddings as soft prompts for input

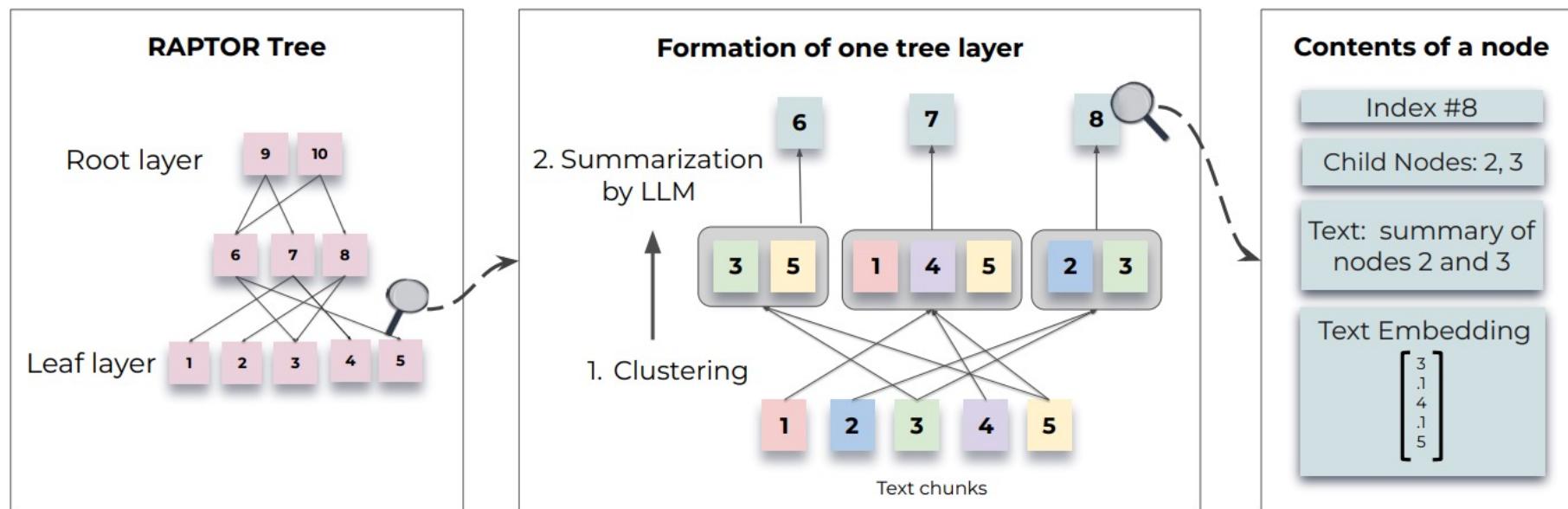


# Level 1 – Explicit Facts - RAG

## Data Processing Enhancement

### Chunking Optimization

- Fixed size chunking, recursive chunking, sliding window chunking, paragraph-based chunking, semantic chunking, etc.
- Refine the text into smaller segments that maintain a high degree of information completeness



# Level 1 – Explicit Facts - RAG

---

## Data Retrieval Enhancement - Indexing

- Sparse Retrieval
- Dense Retrieval
- Other Methods

### Sparse Retrieval

- TF-IDF and BM25
- Word matching methods
- Similarity-based matching (KNN)

# Level 1 – Explicit Facts - RAG

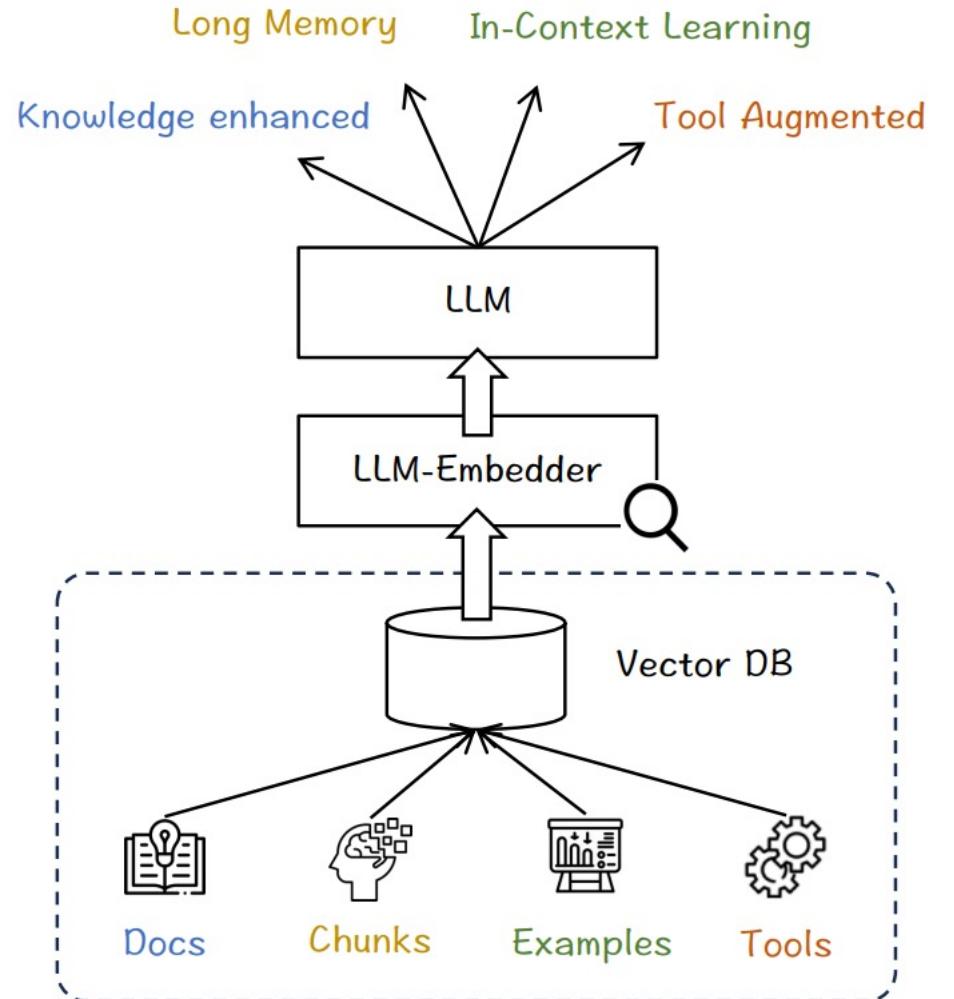
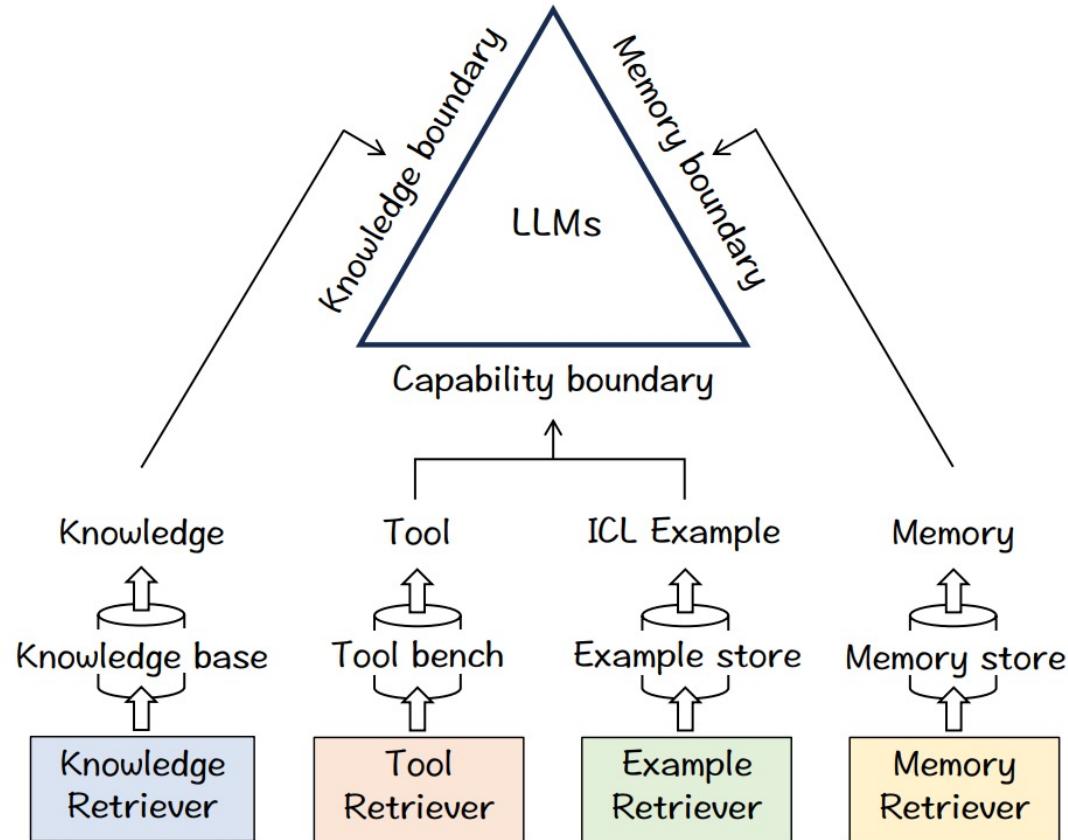
---

## Data Retrieval Enhancement - Indexing

### Dense Retrieval

- BERT-based encoders
- unsupervised contrastive learning for fine-tuning
- feedback from LLMs to guide the training objectives of retrievers
- LLM-based dense retrieval
  - LLM2Vec
  - Llama2Vec

# Level 1 – Explicit Facts - RAG

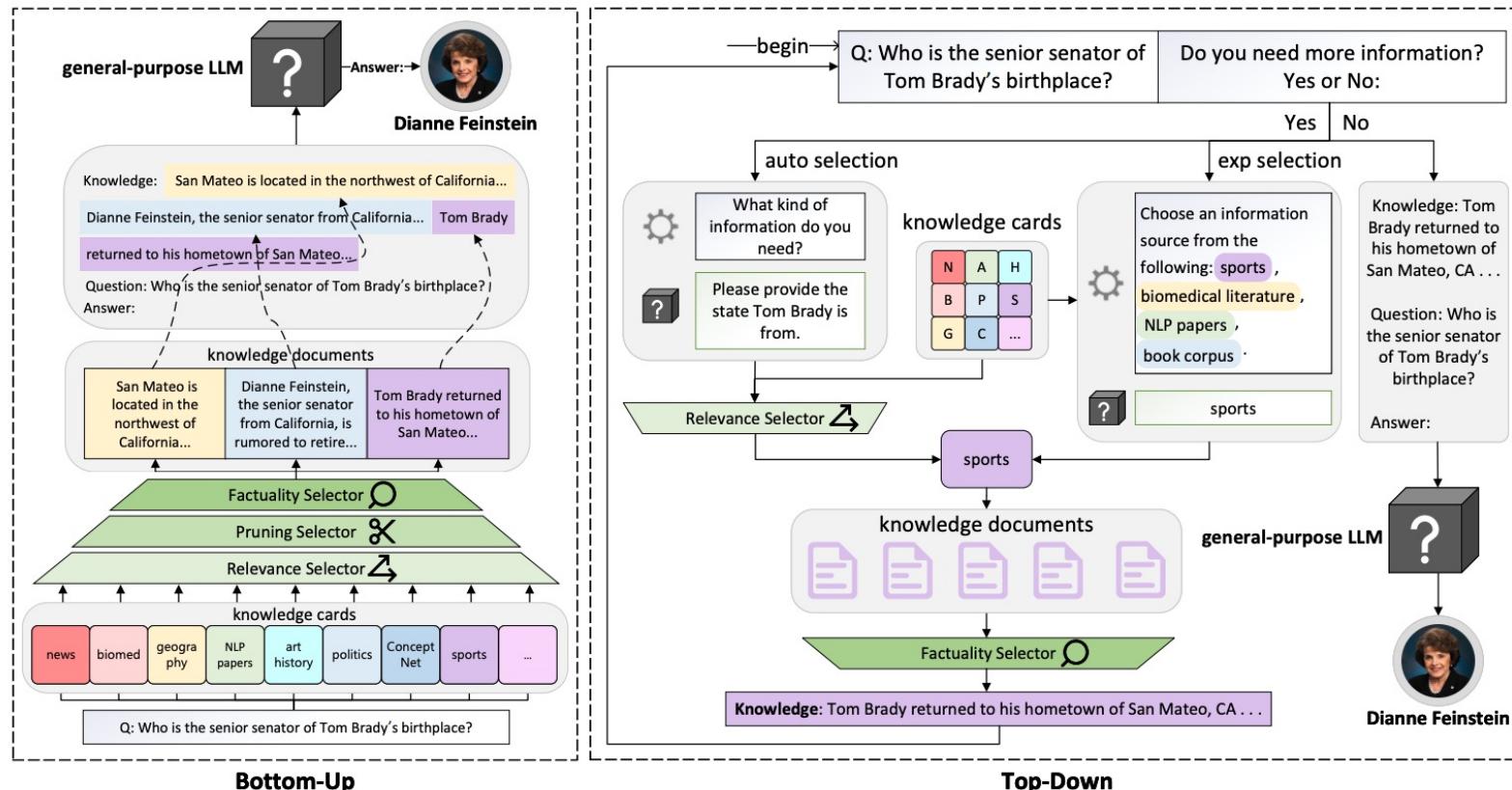


# Level 1 – Explicit Facts - RAG

## Data Retrieval Enhancement - Indexing

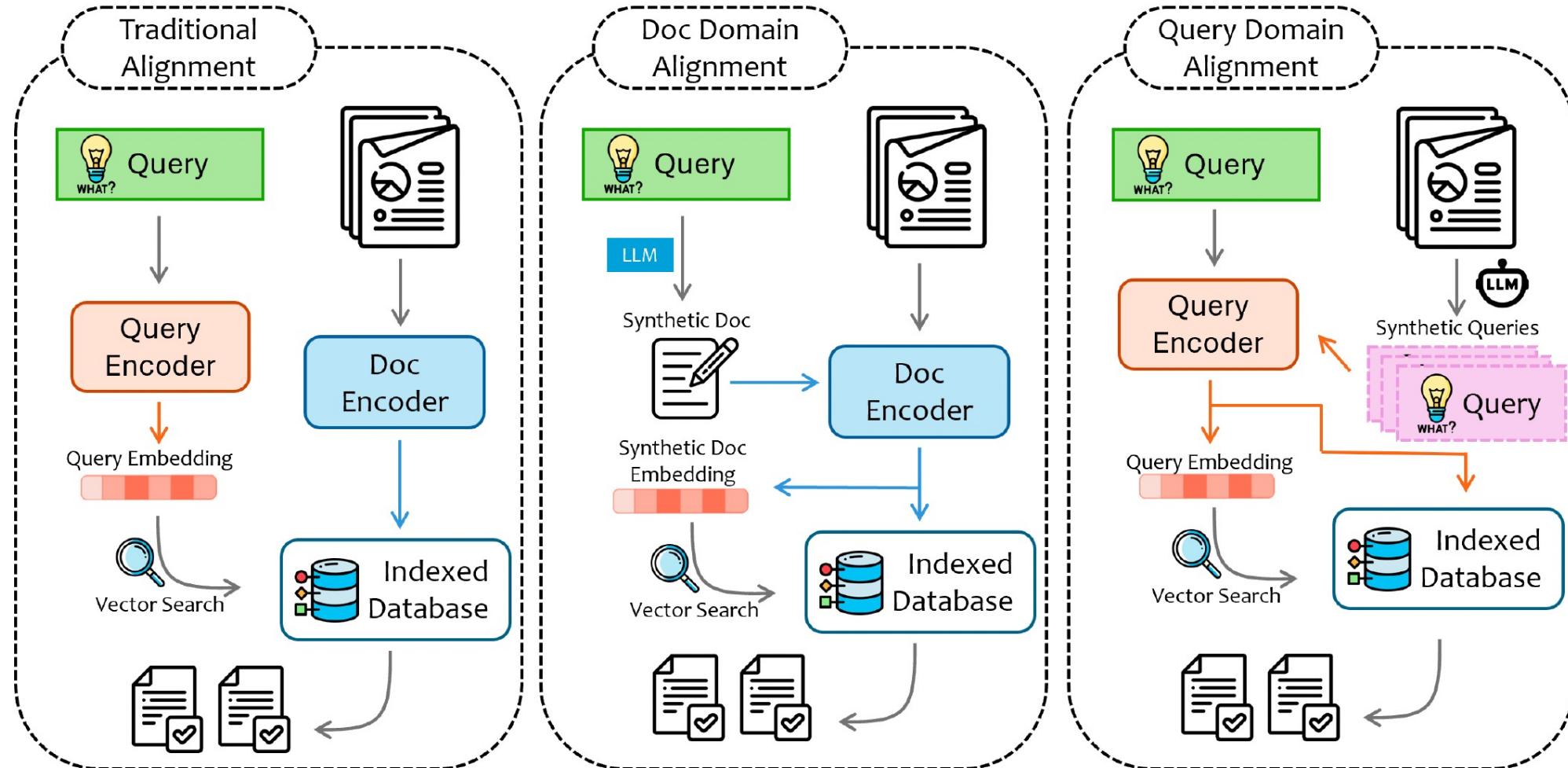
### Others

- Initially determining the knowledge domain needed to answer a query as a fixed area of expertise, and then using dense retrieval to recall supplementary information within this domain



# Level 1 – Explicit Facts - RAG

## Data Retrieval Enhancement – Query Document Alignment

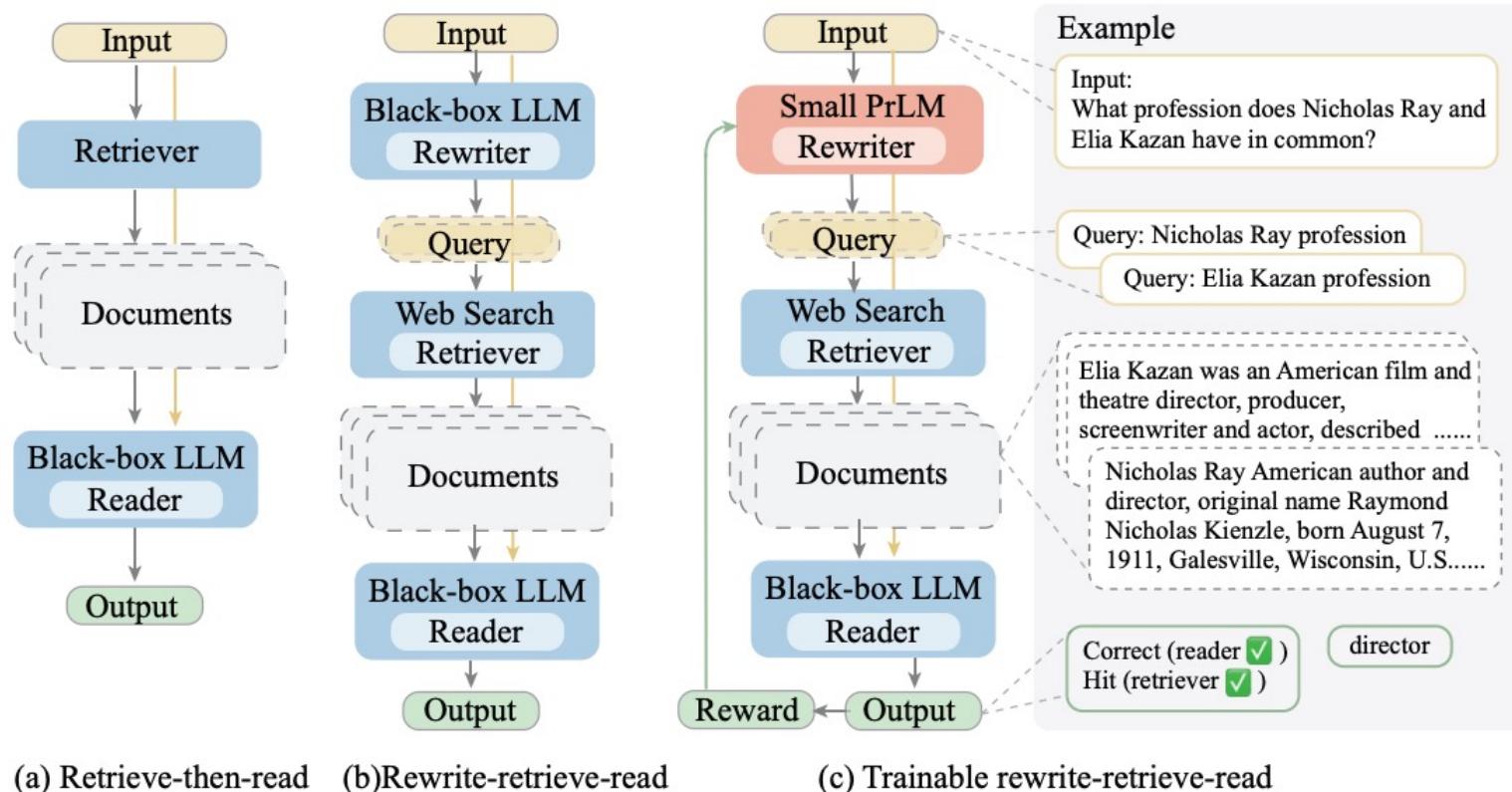


# Level 1 – Explicit Facts - RAG

## Data Retrieval Enhancement – Query Document Alignment

### Traditional Alignment

- Mapping both document segments and the query into the same encoding space.

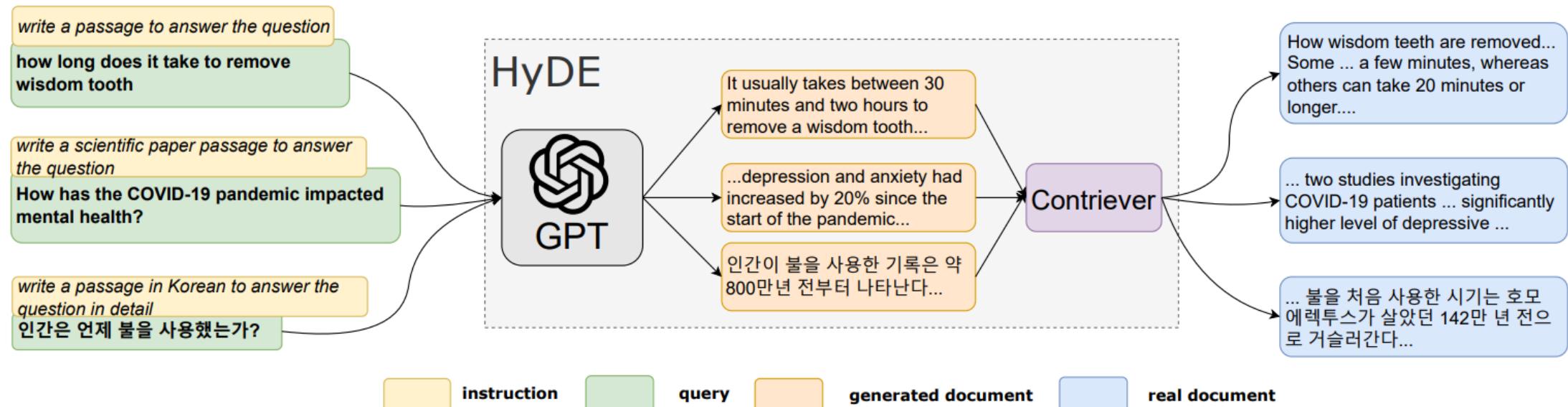


# Level 1 – Explicit Facts - RAG

## Data Retrieval Enhancement – Query Document Alignment

### Document Domain Alignment

- generating synthetic answers first, then using these answers to recall relevant data

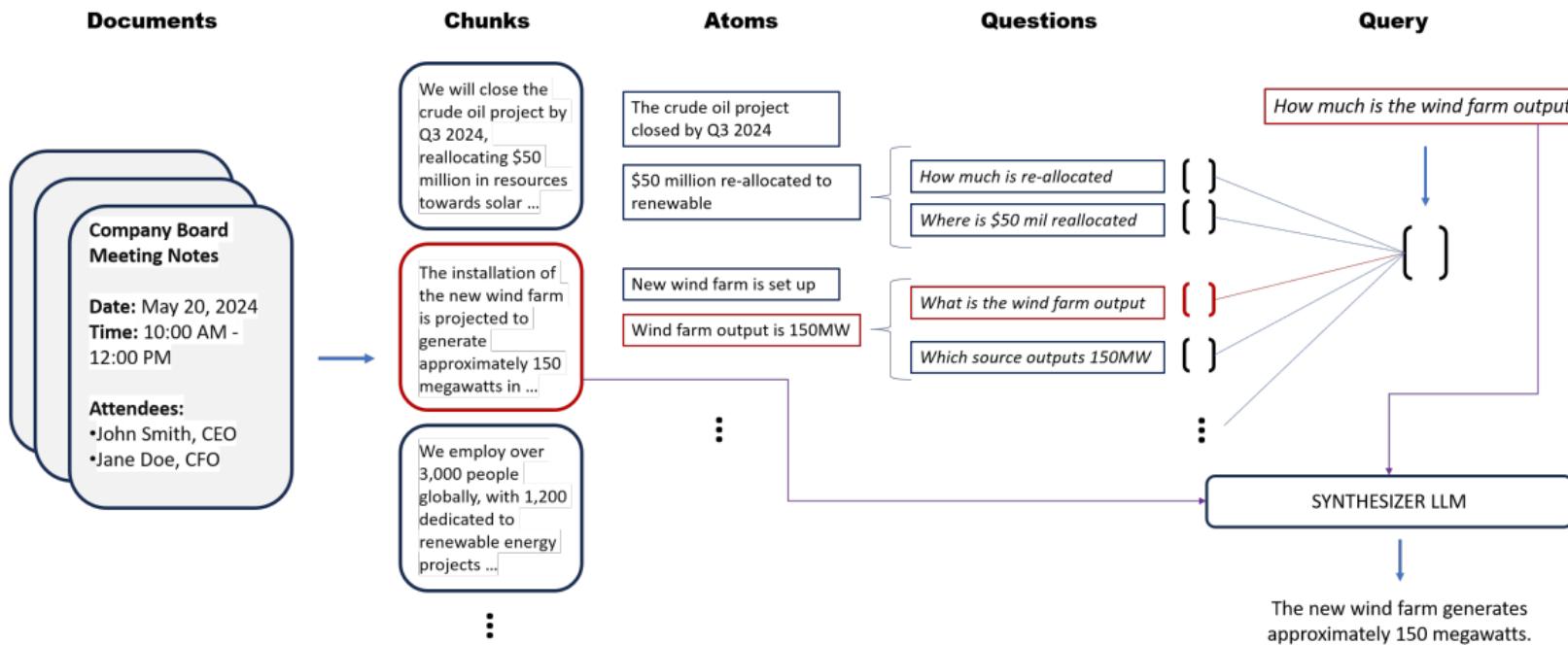


# Level 1 – Explicit Facts - RAG

## Data Retrieval Enhancement – Query Document Alignment

### Query Domain Alignment

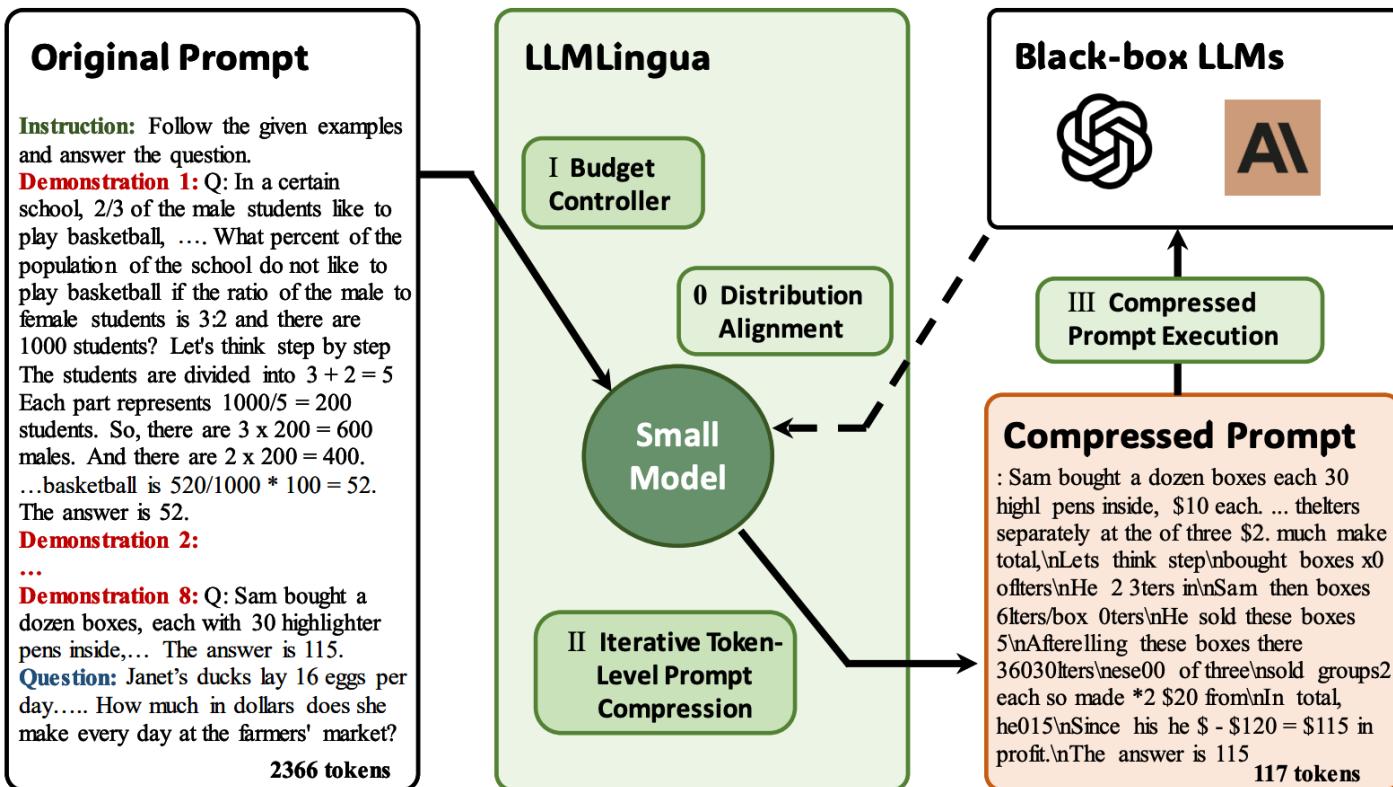
- generating a set of synthetic questions for each atomic unit of text,
- mapping text segments into the query space,
- retrieving the synthetic questions closest to the original query along with their corresponding text segments.



# Level 1 – Explicit Facts - RAG

## Data Retrieval Enhancement – Re-ranking & Correction

- Filter and reorder the segments
- Relevance scores for re-ranking
- perplexity or perplexity gain as ranking criteria



# Level 1 – Explicit Facts - RAG

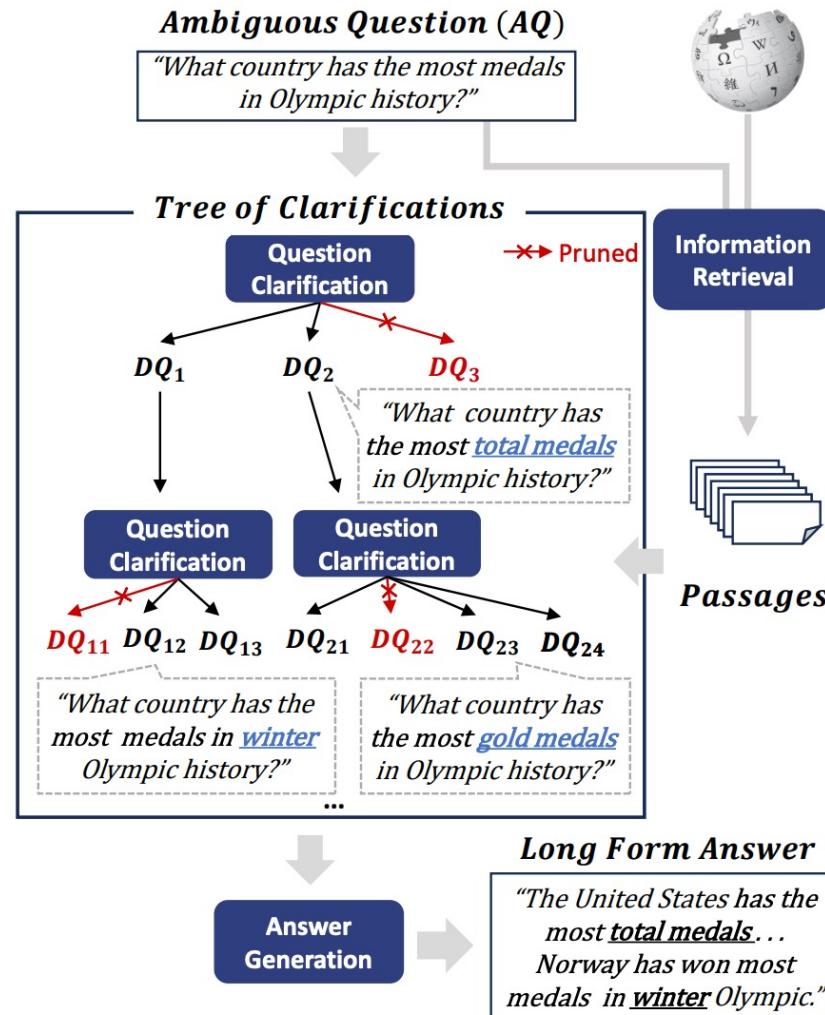
---

## Data Retrieval Enhancement – Recursive Retrieval of Iterative Retrieval

- Overcome the drawbacks of single retrieval
- Tree-like recursive retrieval method
  - Incorporating pruning strategies to incrementally break down ambiguous questions
  - ultimately arriving at the closest correct answer
- SEATER
  - K-means algorithm to construct a hierarchical tree structure of items

# Level 1 – Explicit Facts - RAG

## Data Retrieval Enhancement – Recursive Retrieval of Iterative Retrieval



# Level 1 – Explicit Facts - RAG

---

## Response Generation Enhancement

- Determining if the retrieved information is sufficient or if additional external data is needed
- Handling conflicts between retrieved knowledge and the model's internal prior knowledge
- Supervised Fine-Tuning
- Training data with irrelevant retrieval noise, relevant retrieval noise, and counterfactual retrieval noise
- Joint training of both retriever and generator

# Level 2 – Implicit Facts

---

## Overview.

- Requires common sense reasoning or basic logical deductions
- Information from multiple segments or require simple inferencing

"What is the majority party now in the country where Canberra is located?"

- Decomposing the original query into multiple retrieval operations
- The aggregation of results

## Example.

- Statistical queries, descriptive analysis queries, and basic aggregation queries
- Counting, comparison, trend analysis, and selective summarization are common in "how many" and "what's the most" type queries (multi-hop queries)

# Level 2 – Implicit Facts

---

## Challenges and Solutions.

### Adaptive Retrieval Volumes

- Questions require varying numbers of retrieved contexts
- Specific number of retrieved contexts can depend on both the question and the dataset
- Fixed number of retrievals may result in either information noise or insufficient information

### Coordination between Reasoning and Retrieval

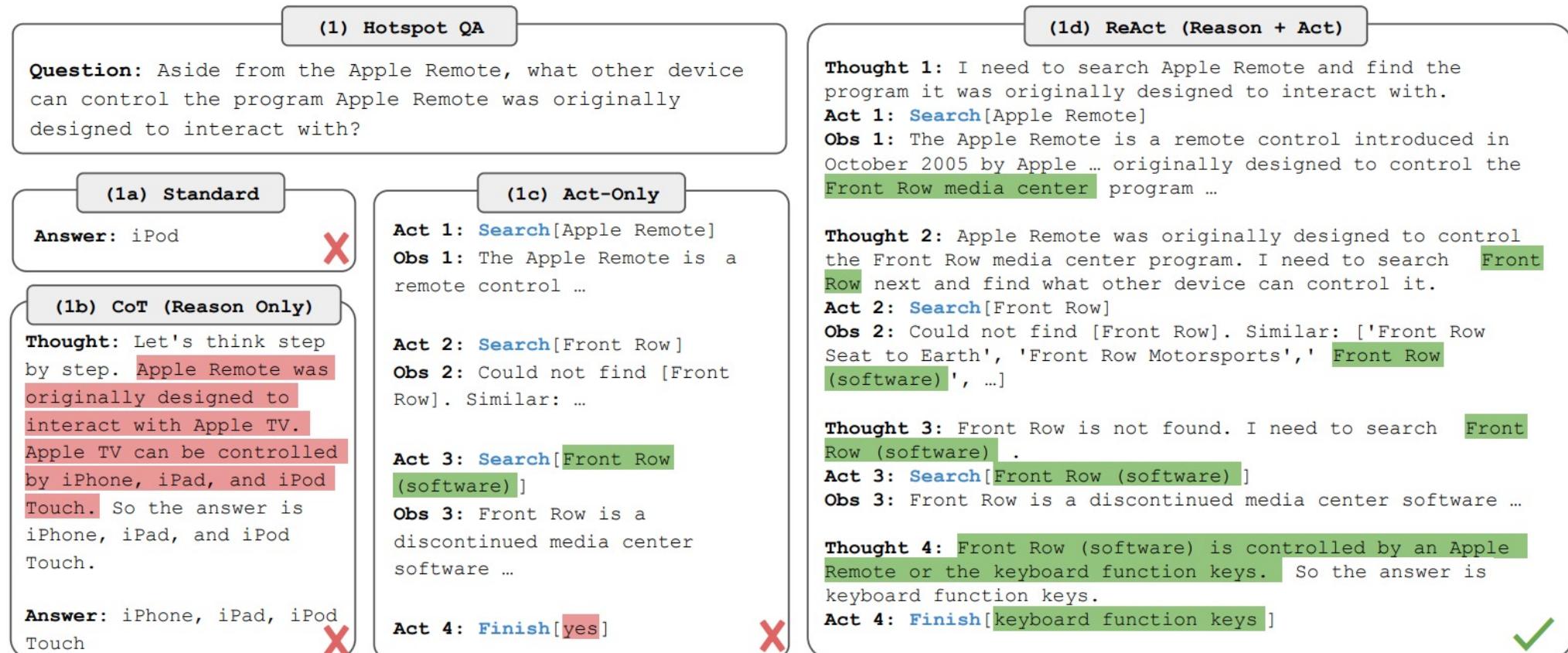
- Reasoning can guide retrieval
- Insights from retrieved information can iteratively refine reasoning strategies

# Level 2 – Implicit Facts

## Iterative RAG.

### Planning Based

- Generating stepwise retrieval plan during the prior retrieval stage or dynamically within the retrieval process



# Level 2 – Implicit Facts

## Iterative RAG.

### Information Gap Filling Based

Question:  $q$

What is the height of the player who won the 2015 AFL Rising Star award?

Retrieval:  $q \rightarrow \mathcal{D}_q$

Title: 2015 AFL Rising Star Context: The NAB AFL Rising Star award is given annually to a stand out young player in the Australian Football League. The award was won by Jesse Hogan of Melbourne

Retrieval-Augmented Generation:  $\mathcal{D}_q \parallel q \rightarrow y_1$

The 2015 AFL Rising Star award was won by Jesse Hogan of Melbourne. Jesse Hogan is a professional Australian rules footballer. He is 198 cm tall. So the answer is 198 cm

Iteration 1

Question:  $q$

What is the height of the player who won the 2015 AFL Rising Star award?

Generation-Augmented Retrieval:  $y_1 \parallel q \rightarrow \mathcal{D}_{y_1 \parallel q}$

Title: Jesse Hogan Context: Jesse Hogan ... playing for the Melbourne Football Club. A key forward, Hogan is 1.95 m tall ... made his AFL debut in the 2015 season and won the Ron Evans Medal as the AFL Rising Star

Retrieval-Augmented Generation:  $\mathcal{D}_{y_1 \parallel q} \parallel q \rightarrow y_2$

The 2015 AFL Rising Star award was won by Jesse Hogan of Melbourne. Jesse Hogan is 1.95 m tall. So the answer is 1.95 m

Iteration 2

# Level 3 – Interpretable Rationales

---

## Challenges and Solutions.

### Prompt Optimization Cost

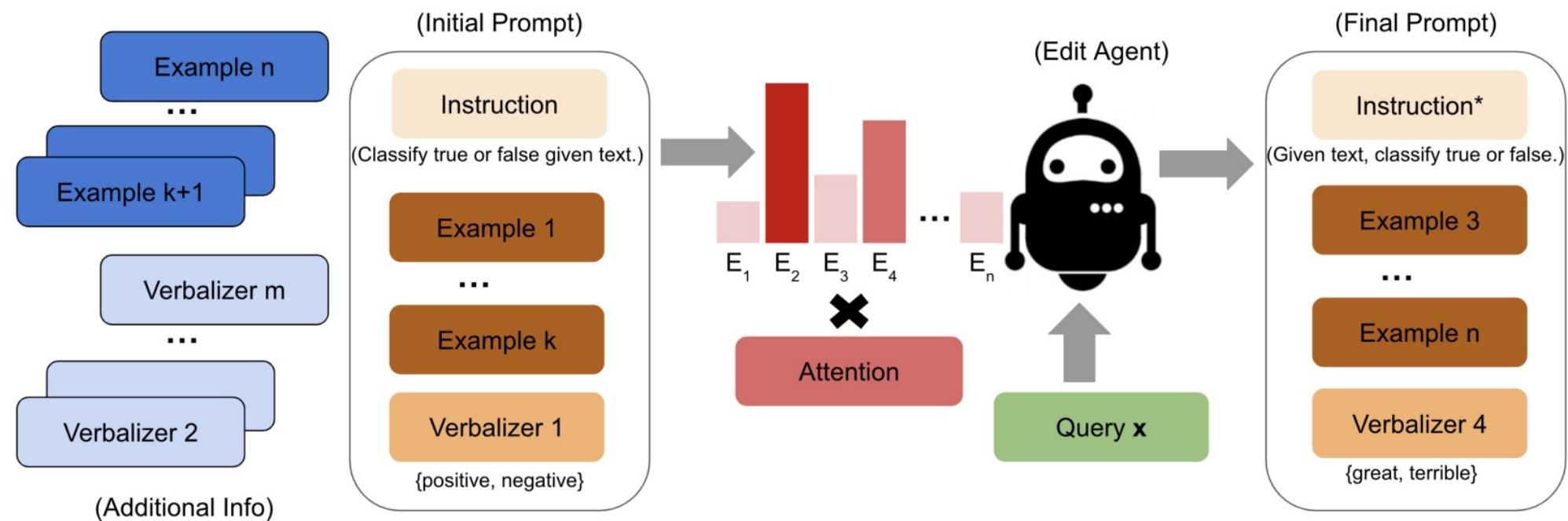
- Queries demand tailored background knowledge and decision-making criteria
- Manually designed prompts can be highly effective but are labor-intensive and time consuming
- Training models to generate tailored prompts incurs significant computational overhead

### Limited Interpretability

- Impact of prompts on LLMs is opaque
- Access to the internal parameters of LLMs is restricted

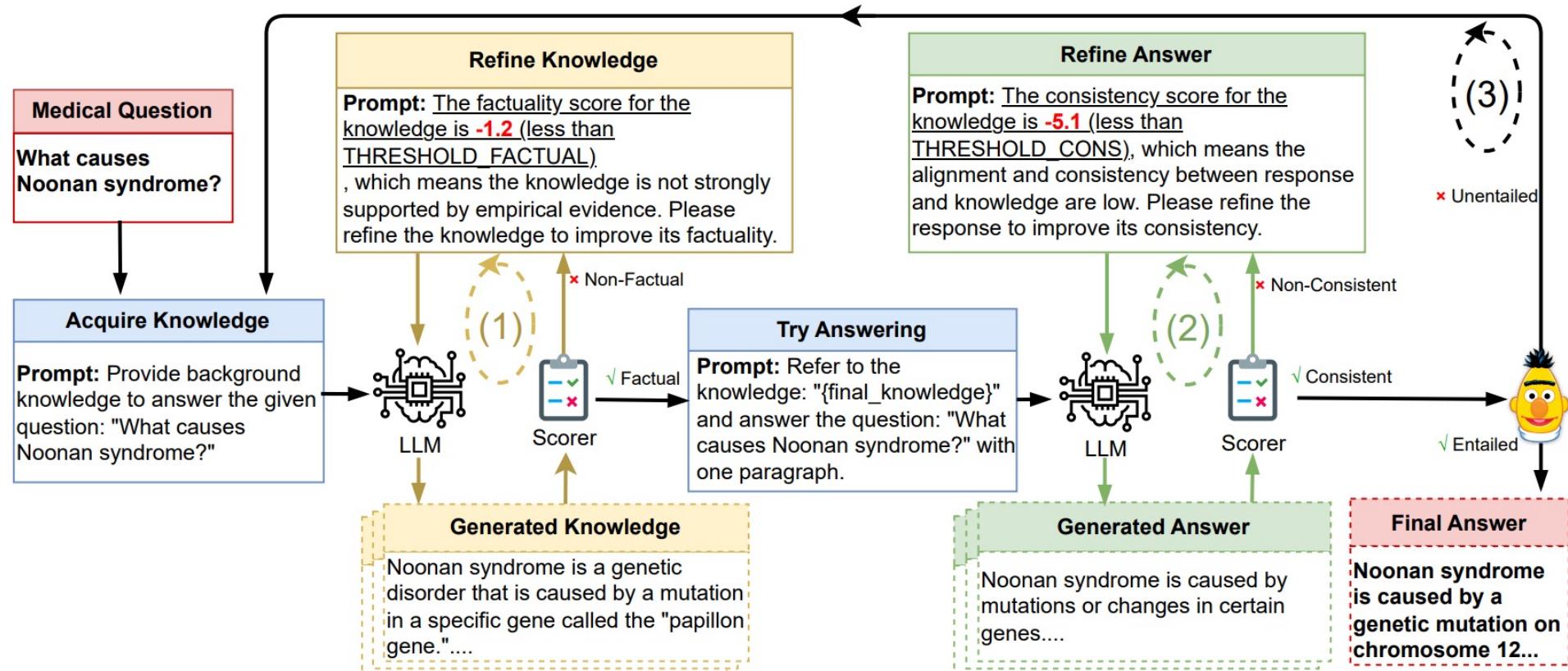
# Level 3 – Interpretable Rationales

## Prompt Tuning.



# Level 3 – Interpretable Rationales

## Chain of Thought.



# Level 4 – Hidden Rationales

---

## Overview.

### In-domain Data

- utilize data from the same domain
- historical question-and-answer records or artificially generated data.
- contains the reasoning skills or methodologies
- e.g., Python programming puzzles

### Preliminary Knowledge

- extensive, dispersed knowledge bases
- comprehensive axiomatic system
- e.g., all local legal codes that form the basis for legal judgments
- intermediate conclusions that simplify reasoning processes
- e.g., mathematical proofs

# Level 4 – Hidden Rationales

---

## Challenges and Solutions.

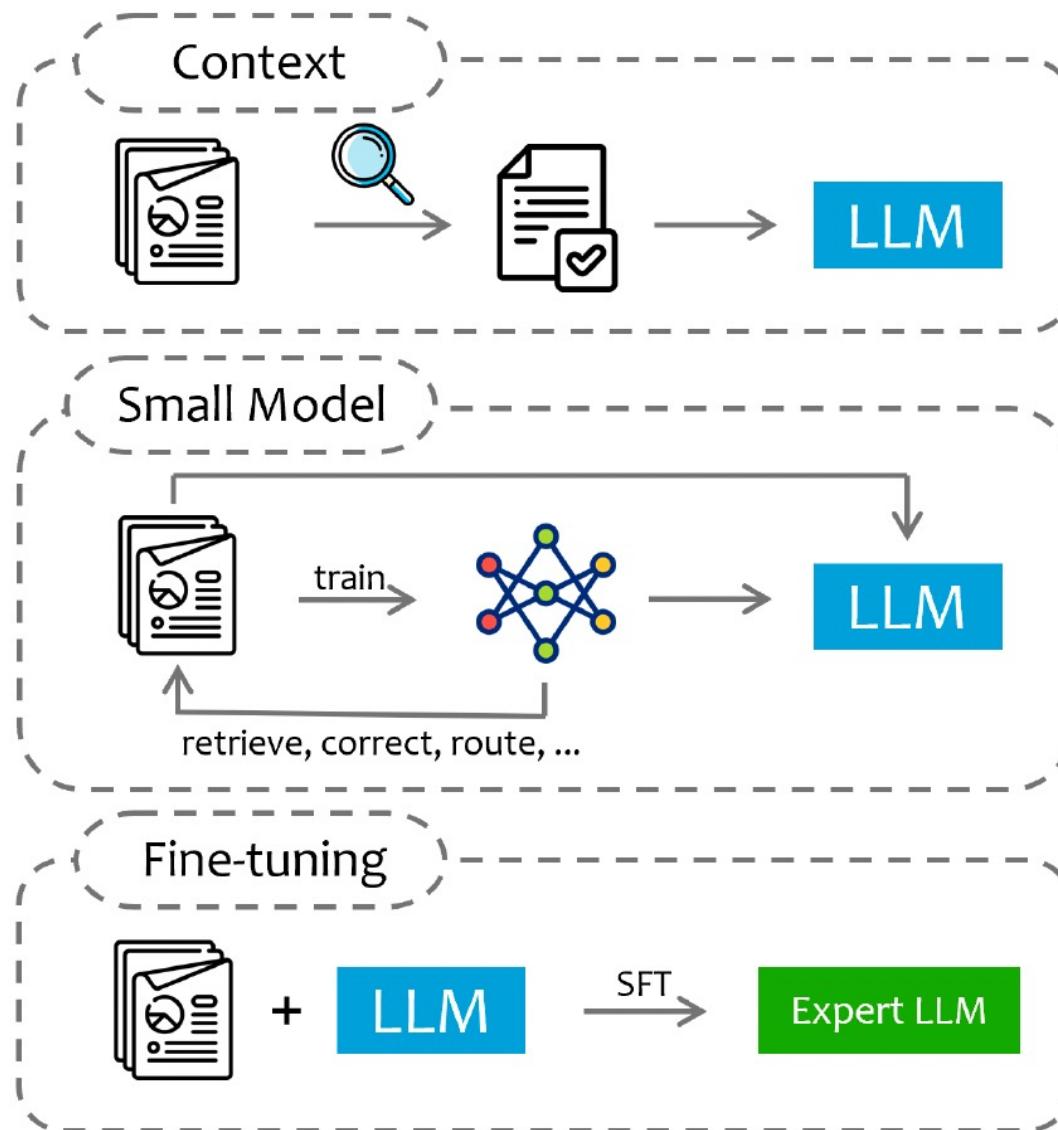
### Logical Retrieval

- Struggle to capture the true target of the query
- Need for parse and identify underlying logical structures

### Data Insufficiency

- Relevant information is often dispersed knowledges
- Demands robust capabilities in data interpretation and synthesis

# Level 4 – Hidden Rationales



# Evaluating RAG Systems

---

## Metrics for Evaluating RAG Systems

- Context Precision: the proportion of relevant **chunks** in the retrieved contexts
- Context Recall: how many of the **relevant documents** were successfully retrieved
- Noise Sensitivity: how often a system makes **errors** by providing incorrect responses when **utilizing** either relevant or irrelevant retrieved documents
- Response Relevancy: assessing how pertinent the generated answer is to the **given prompt**
- Faithfulness: factual consistency of the generated answer against the **given context**