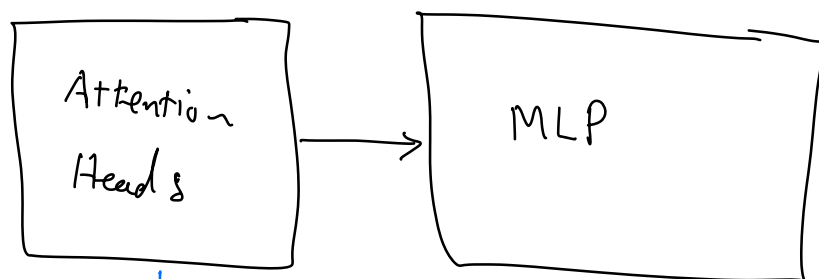
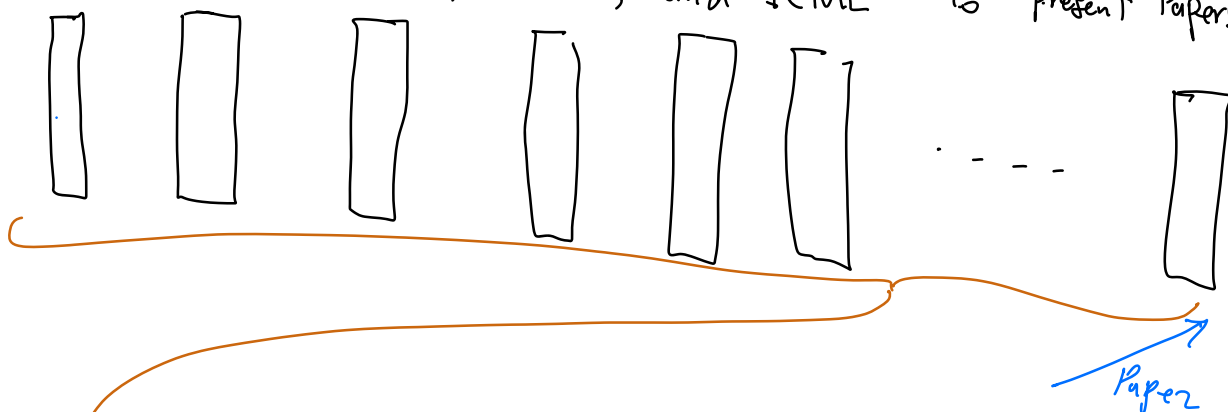
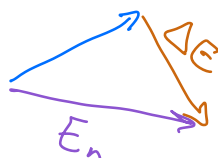


I attended KDD, VLDB, and ICML to present Papers.



ΔE



Databases
↓
Non-Parametrized
Retrieval-Only
Systems

Retrieval-Augmented
Generation
(RAG)

LLMs
↓
Data Generation
only based on the
Learned Parameters

Parametrized

- Always returns facts ✓

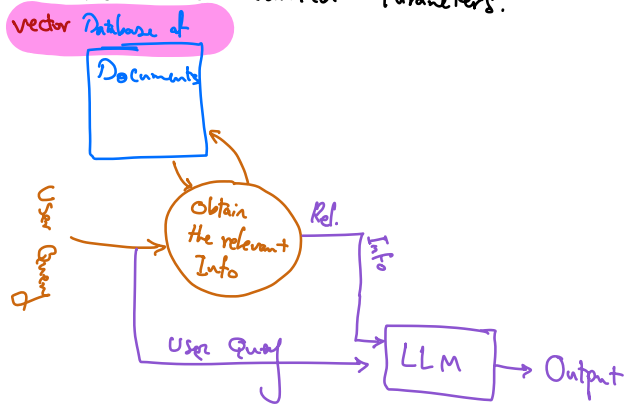
- Does not use the Collective
Information In The world

RAG

- May return wrong Info.

- Generated based on
Rich parameters learned
by the model ✓

RAG: Separates the factual Statements from the Learned Parameters.



Challenges:

- Retrieval Quality

- High Precision: The retrieved documents are highly relevant to the Query

+ Fine-Tune Embeddings based on the context.

+ Use an Ensemble of Fine-Tuned Embedders

- High Recall: The relevant documents are not missed

+ Optimize the chunking

+ Increase the value

of k

→ will lead to lower Precision

→ A lot of noise (irrelevant Documents) will be retrieved

Challenge 2:

Integrating the retrieved documents into the Prompt.

- The noise in the retrieved documents can mislead the LLM

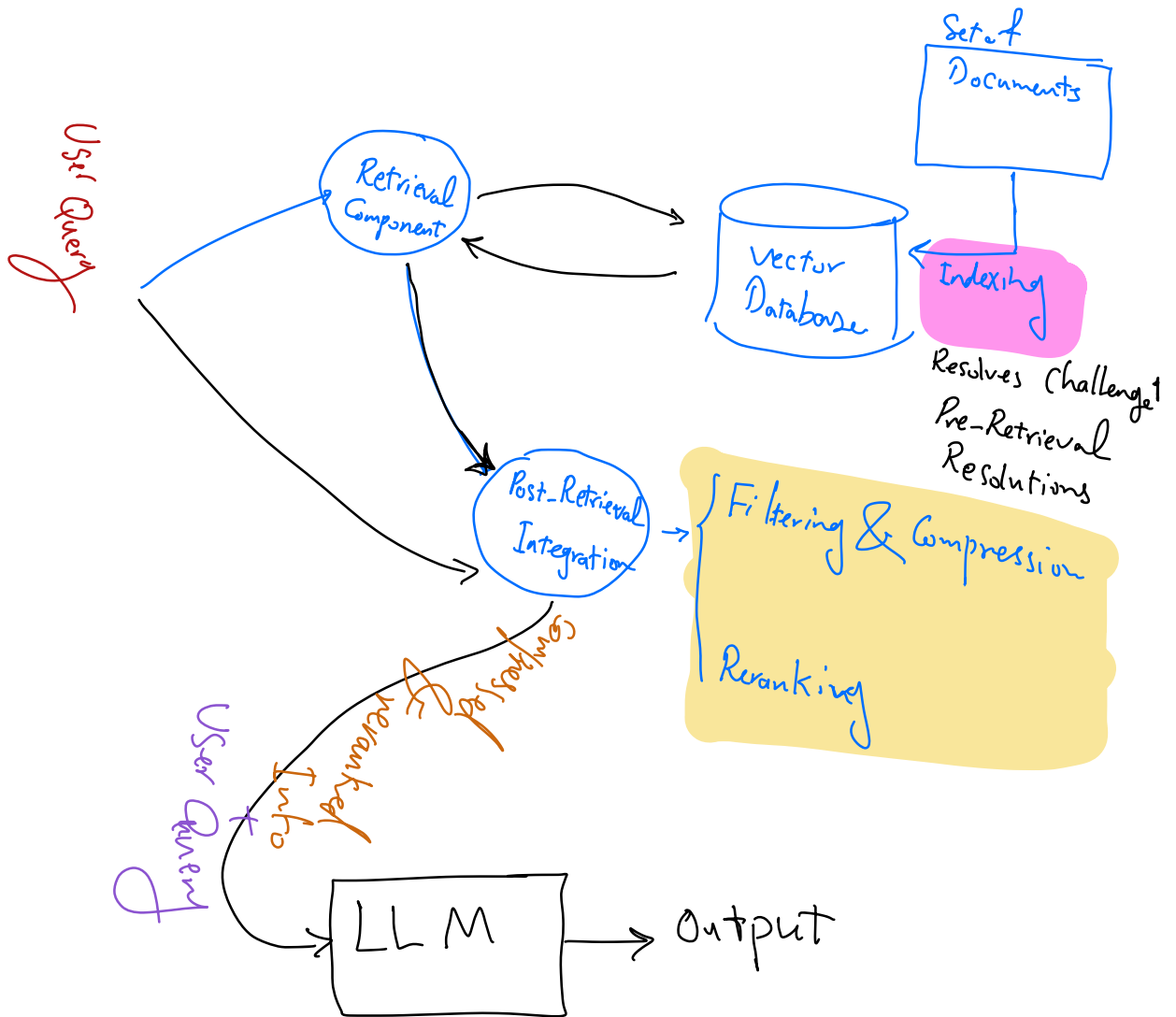
- Limited context window size

- Lost In The Middle issue

Challenge 3:

The LLM may still put a high weight on the Learned Parameters (MLP)

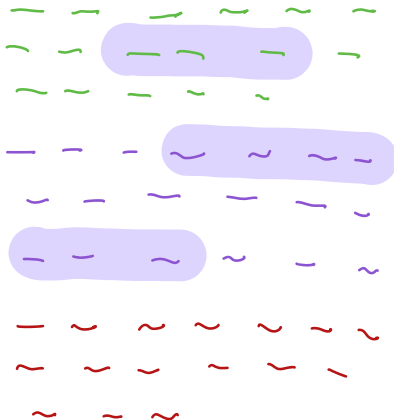
+ Retrain / Fine-Tune the LLM for RAG



Filtering & Compression:

to take the retrieved Sequence of Documents
& remove the noise & extract the relevant
Information

e.g., Retrieved Chunks



Step 1: Filtering



Step 2: Extract Relevant Parts of each chunk

