

RAG [Retrieval Augmented Generation]

Why we need RAG?

- ① Hallucination $\xrightarrow[\text{with}]{\text{dealing}}$ ① temperature
↑
model being over-creative and not factual
- ② Fine-tune a LLM (Creating Distilled models)
(Pefit, LoRA)
- ③ Adding Context
 { - Using Prompt Template (FewShotTemplate)
 - Using memory data
- ② Stale Knowledge
- ② Size Limits
- ④ Specific Context Awareness.

RAG helps us:

- ① Make context RECENT
- ② more accurate and on-point answers.
- ③ Uses your own data to get answers.

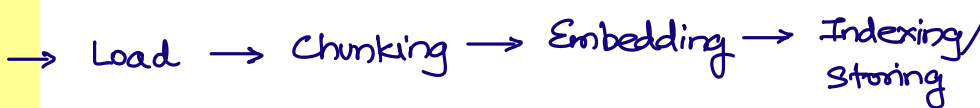
What is RAG?

RAG is a pattern where LLM answers a question using

EXTERNAL UP-TO-DATE PROPRIETARY DOCUMENTS

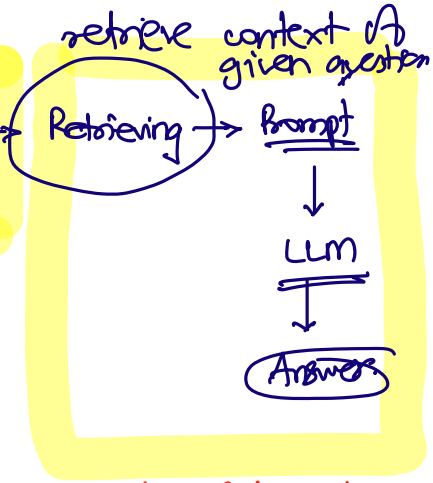
that are retrieved at query-time.

Core Data Source



pdf, web page,
wiki, SQL rows,
etc

Data Indexing



Data Retrieval and Generation

Embeddings

embedDim = 2

Vector form

