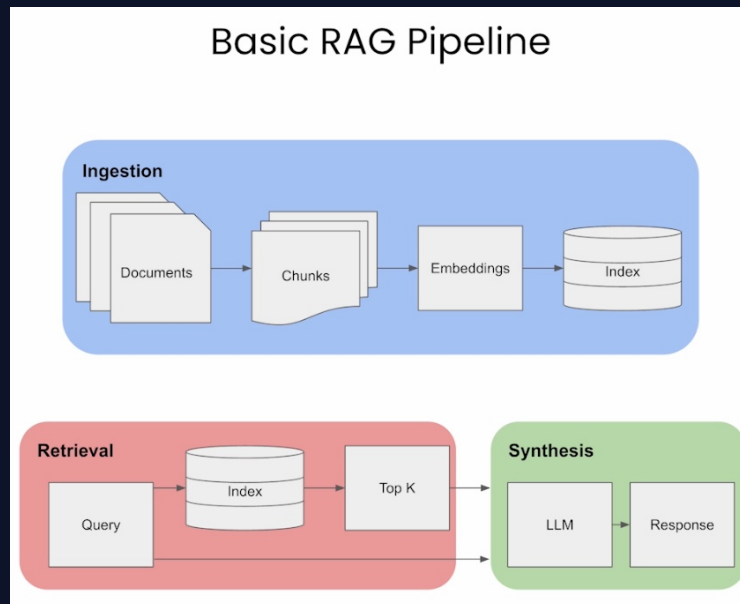


# What is RAG:

## Understanding Retrieval-Augmented Generation



Agenda: What is RAG?

Why?

Importance?

Use case?

Benefit?

1

Core concept? } - 2

Components } 3

Challenges

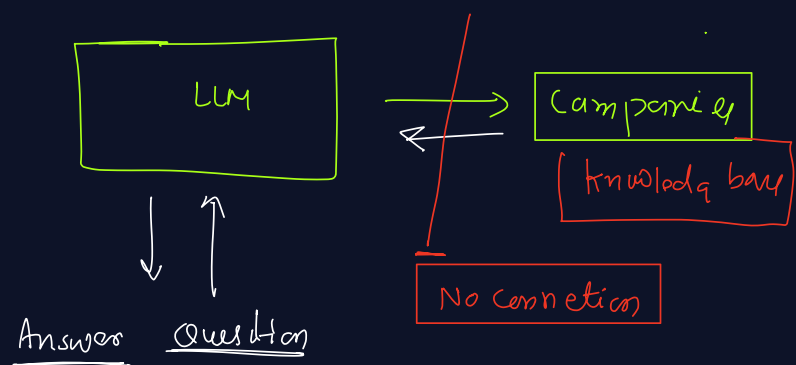
Limitations

Future Trends

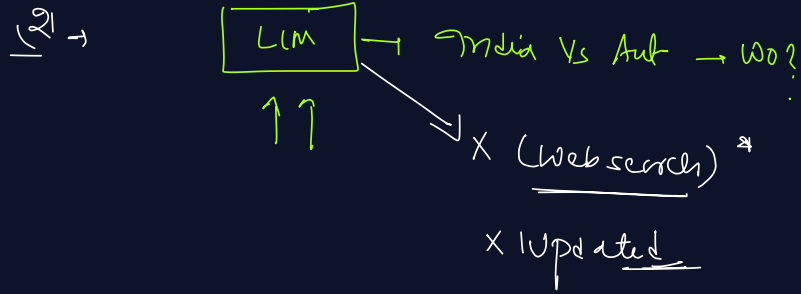
4

Implementation - 5

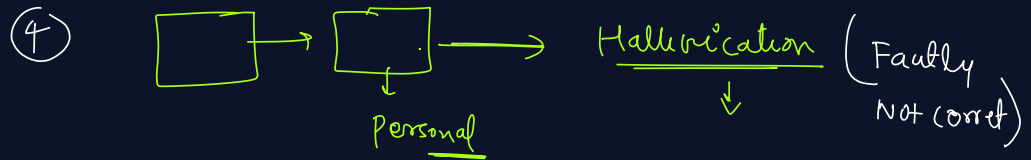
(1), key : 1)



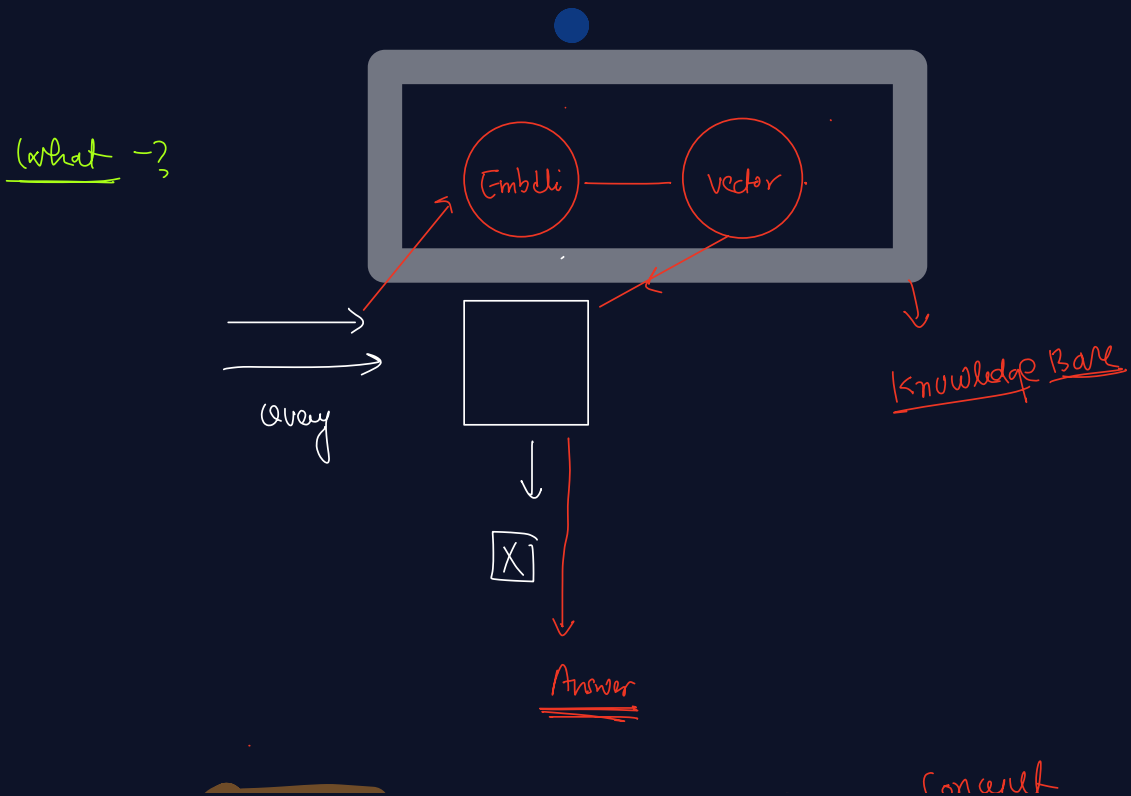
Personal Database

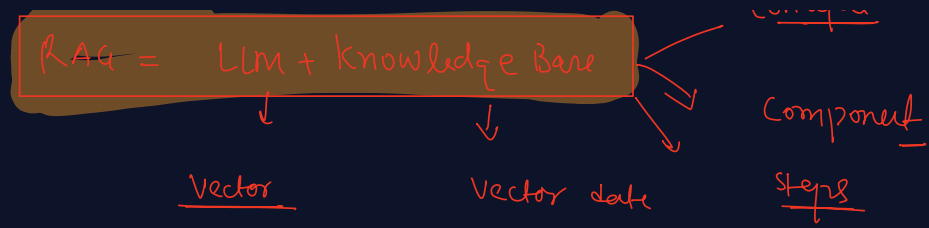


(3) • Real time Data ⇒ Live - commentary X



(5) → Domain specific knowledge





(1) → Benefit & Use Case →

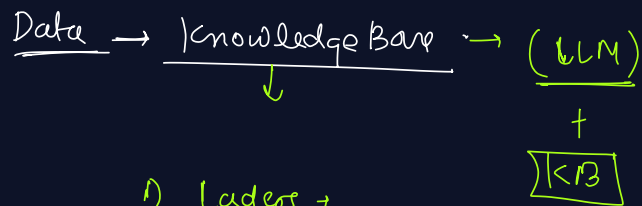
- 1) Domain specific Model
- 2) Real time data Access
- 3) Question - Answering.
- 4) Chats for your Database

RAG ⇒ Concepts:

- 1) Query ⇒ Question
- 2) Retrieval →
- 3) Embedding.
- 3) Vector database
- 4) Data loaders.
- 5) Splitter & chunks
- 6) Augmentation

7) Fusion logic

8) Evaluation Matrix



2) Split, Tokenizer →

3) Embedding → Text → Numerical (Meaning)

4) Vector Data → place to store data

5) Knowledge Base → Retrieval process

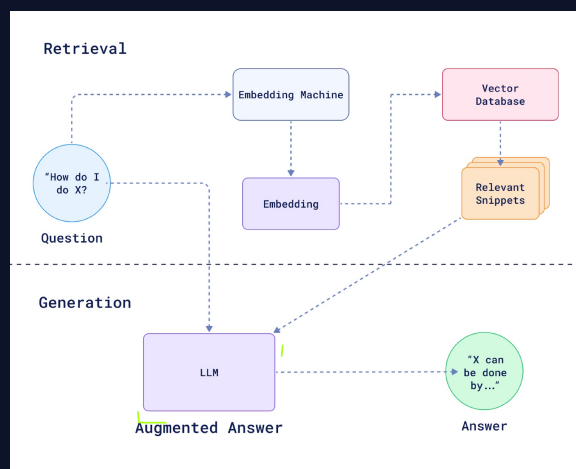
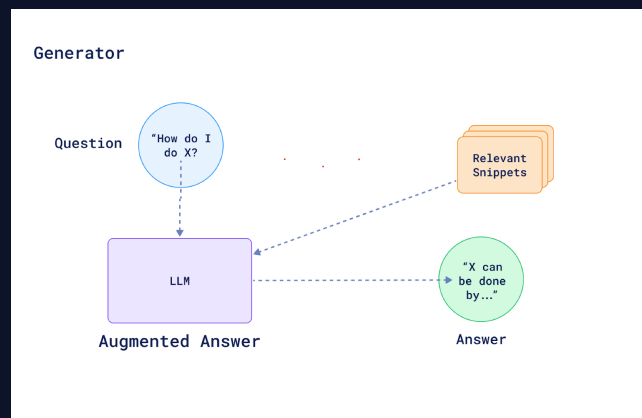
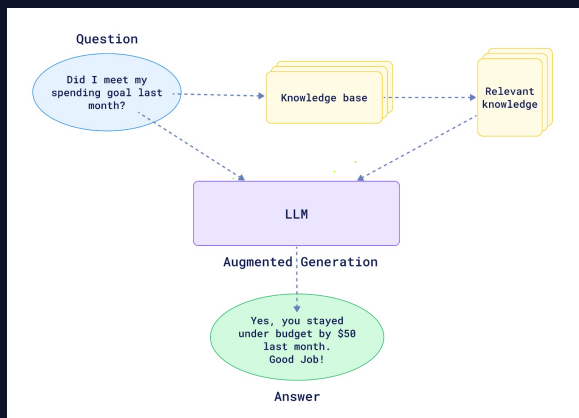
6) Augmented → Combine + Answer.

7) Fusion logic ⇒ Q + KB →

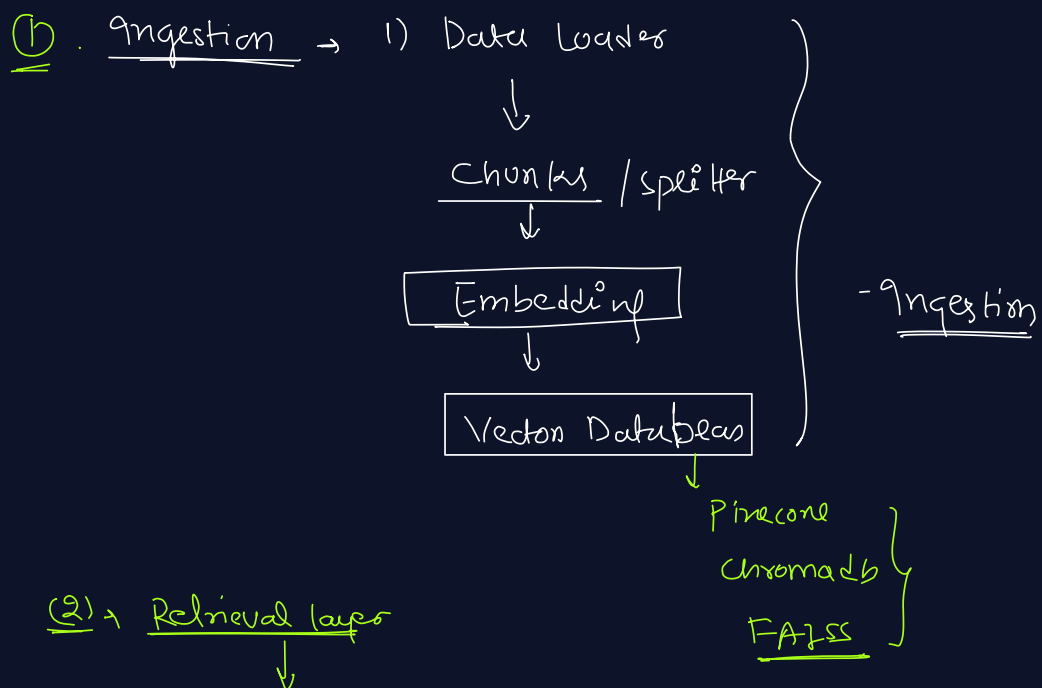
(Q) → Evaluation Matrix ⇒ 1) Answer

- ↓
- 1) Human Answer
  - 2) Confusion
  - 3) RoUGE /
  - 2) BLEU

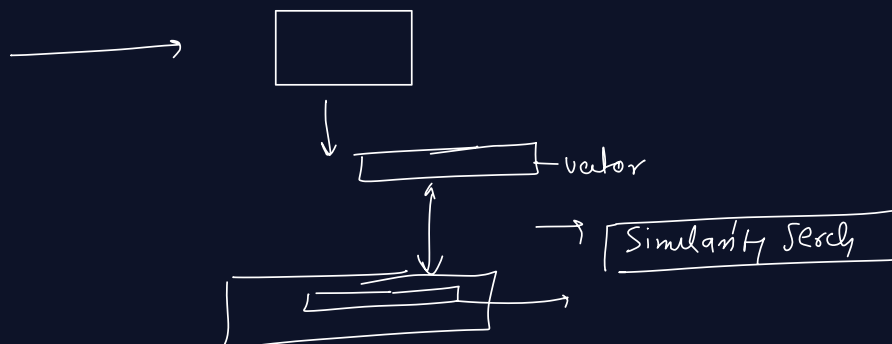
## Architecture & Components

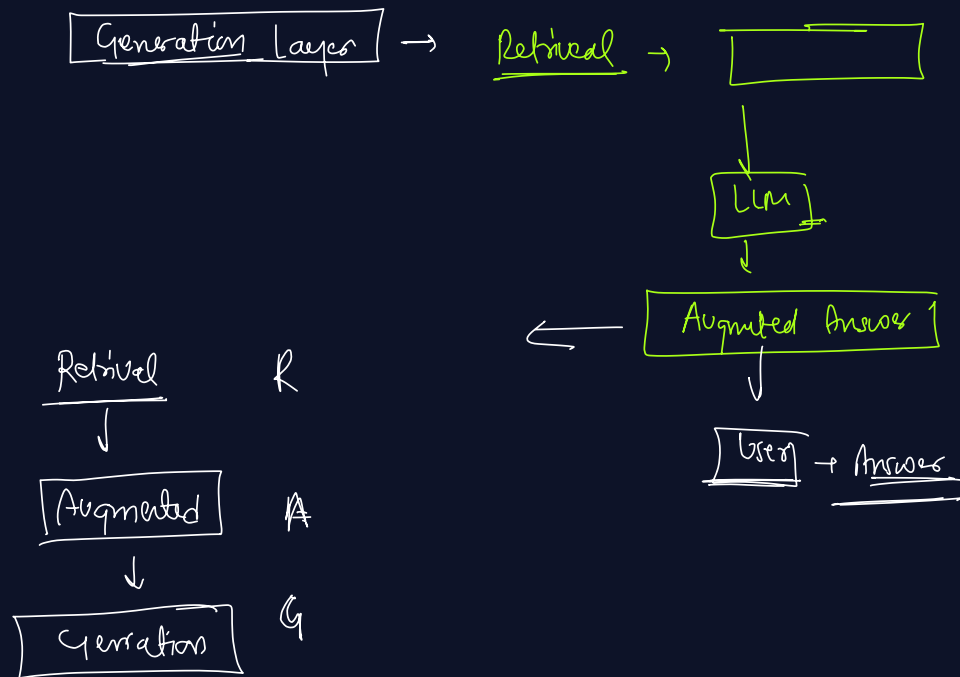


- Components:
- 1) Ingestion →
  - 2) Generation →
  - 3) Retrieval →
  - 4) Data pipeline
  - 5) Interpretation



Cosine Similarity





### Challenges:

1) Dependency of High Quality Data

Knowledge Base

(2) Scalability

(3) Latency (Speed) → Slow (wait to get Respon)

(4) Storage Vectorban → Cost



Expensive

(5) → Argumented Answer → (not correct Answer)