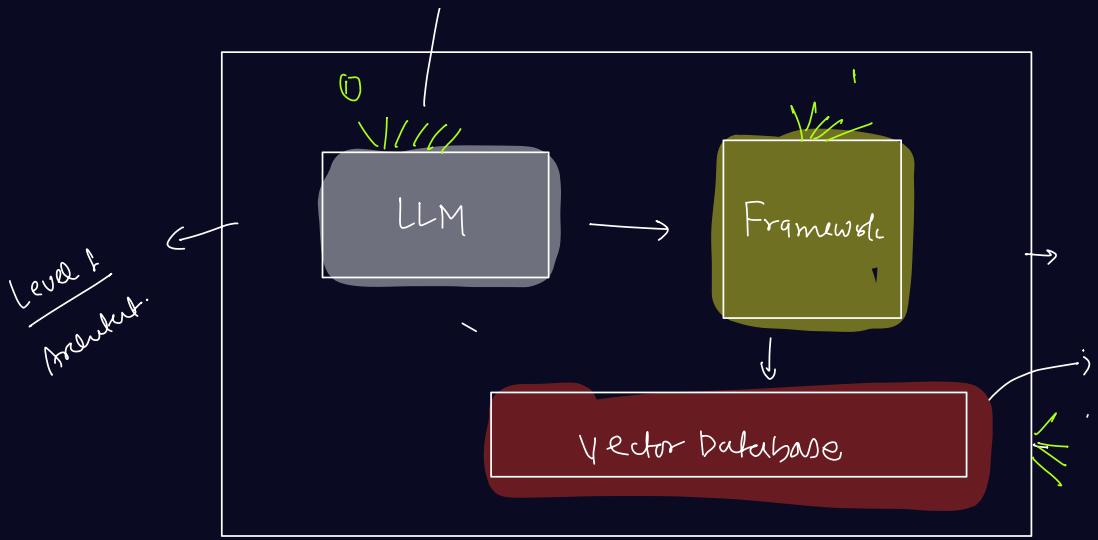
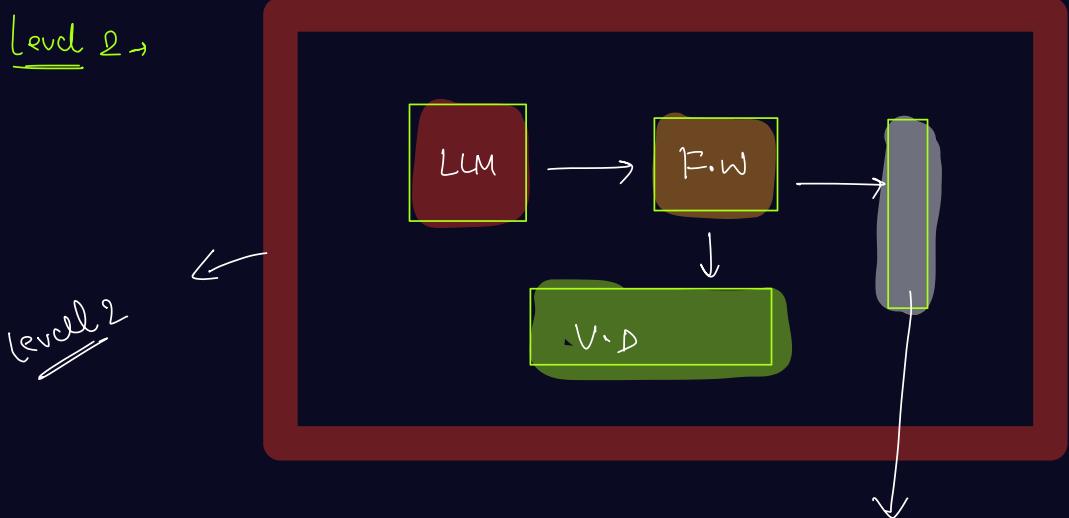


Level 1 architecture → Project → Architecture.

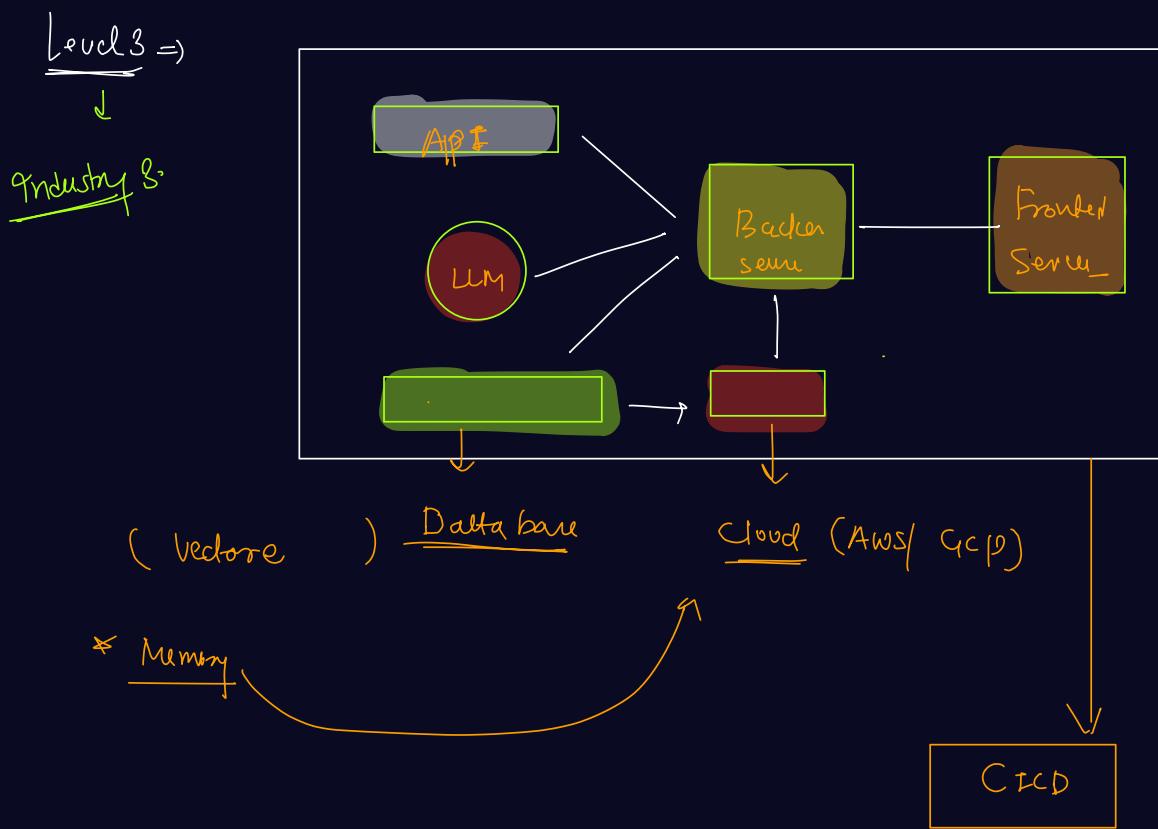
① Langchain

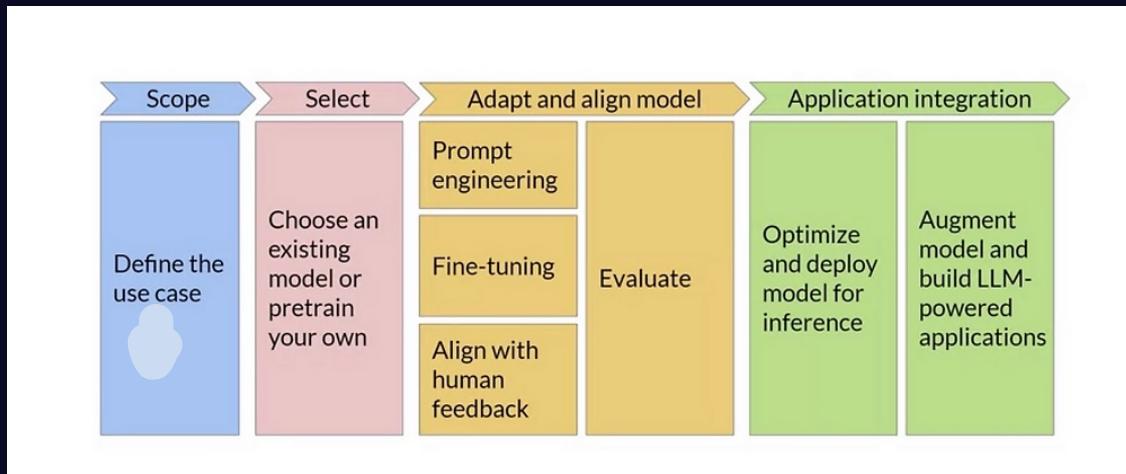


- 1) LLM → Large Language Model
- 2) Orchestration Framework → Langchain / Llamaindex
- 3) Database → Vector Database



- ✓ 1) Flask ← UI / Frontend
- ✓ 2) Streamlit
- ✓ 3) FastAPI





Life Cycle \Rightarrow ①) Model Selection

④ \leftarrow ②) Data Preparing \rightarrow charts
em
vector

③) Technique Generating \rightarrow

① \leftarrow { 1) Prompting
 2) RAG
 3) Finetuning }

4) Evaluation

* Deploy / Integration \rightarrow LLMops

AWS Lambda Vertex

(3) \rightarrow Future Trends \Rightarrow Generative \rightarrow

1) MultiModal \rightarrow Use.

2) AI Agents.

* (3) \rightarrow Responsible AI - (prompting)

(4) \rightarrow Personalised Model \rightarrow

↳ Finetune

↳ RAG

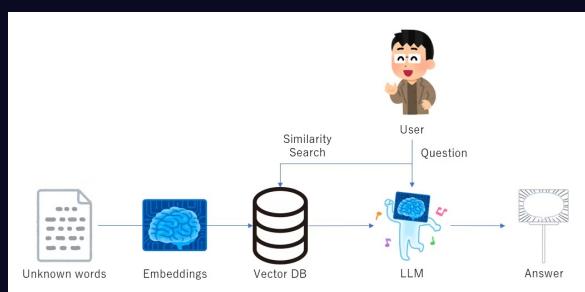
↳ prompting

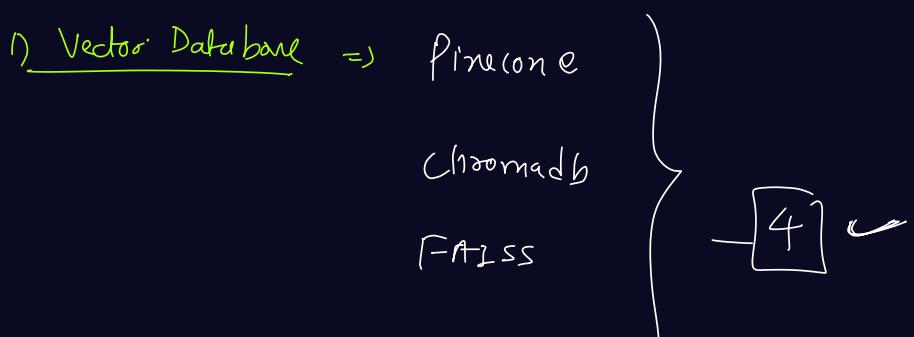
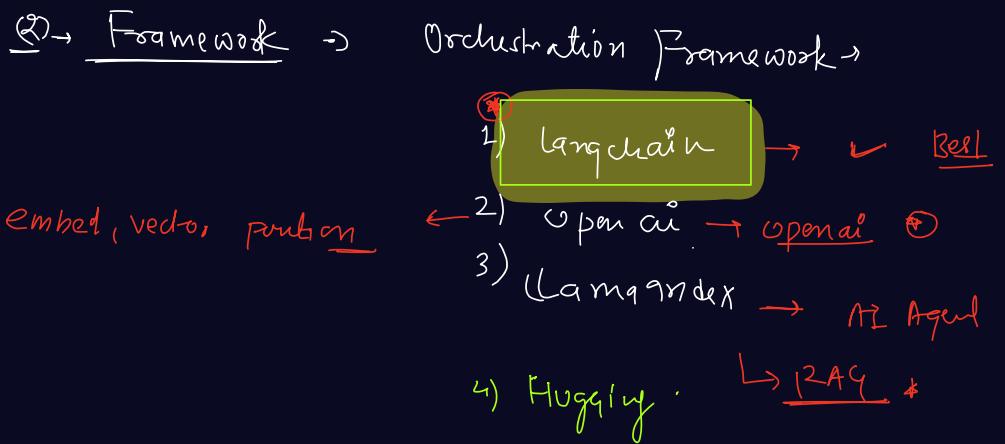
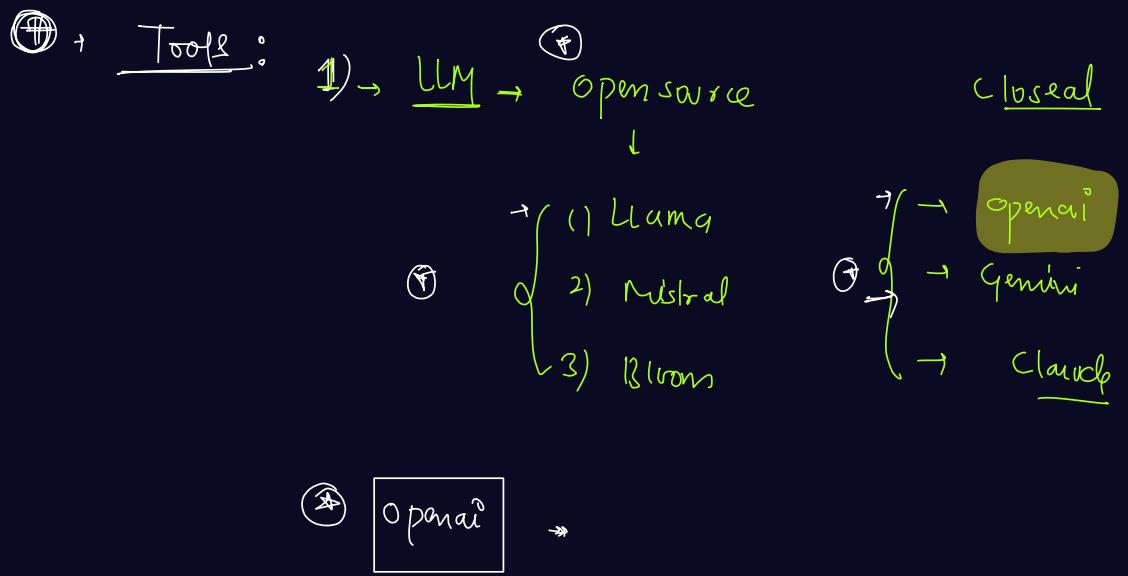
(5) \rightarrow Contextual Window \Rightarrow Memory (

↓

\Rightarrow $\boxed{3}-\boxed{4}$ $\rightarrow M_m$

Buffer Memory / Conversation.





Weational ↴

Level 2 → Frontal Server ↴ 1) Fast API

↳

2) Streamlit

3) Flask

Level 3 → 1) Backend → 1) FastAPI
2) Langchain

2) Cloud Storage ⇒ 1) S3 bucket → AWS

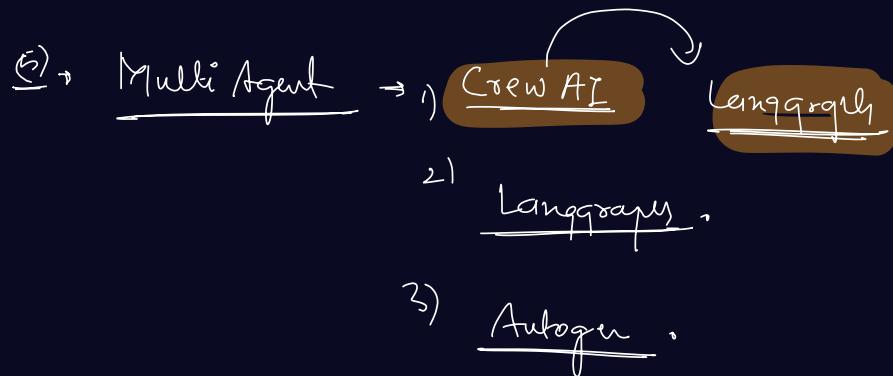
2) GCP → Vertex (Bucket)

(3) Deployment → Github Actions

Circle CI

Jenkins

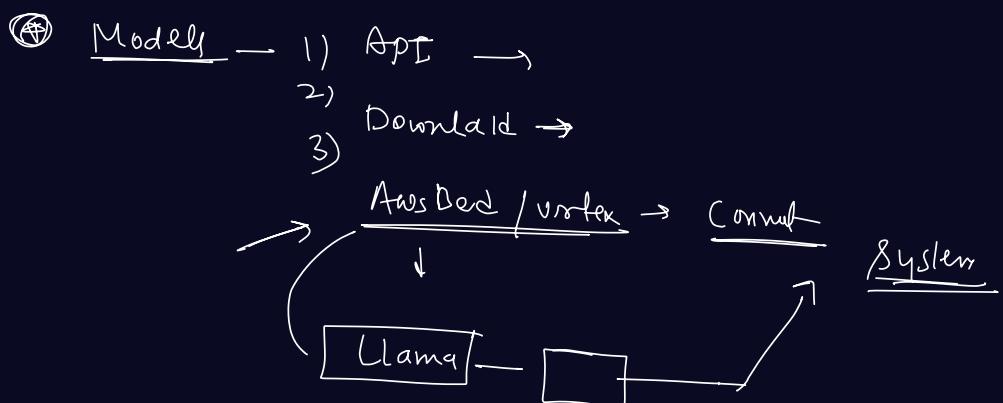
Docker / Container



7 LMMops \Rightarrow (1) Aws-RedRock

Vertex AI

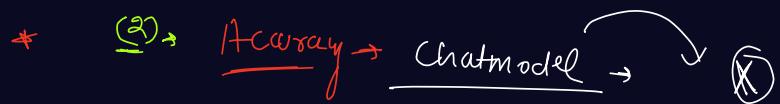
① Lansaphilie



Antivirus

~~Top~~ → How to Select LM Models ?

- ① Requirements → 1) Text
2) Audio → ✅
3) Video
4) Image

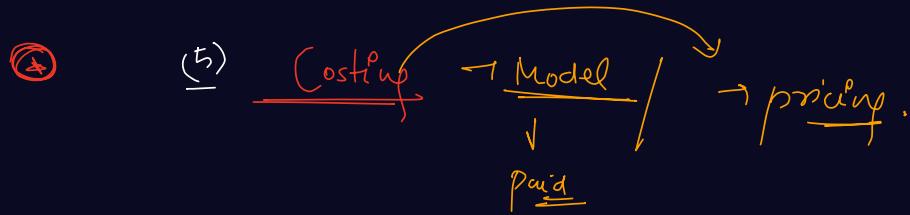


② + $\frac{\text{Latency}}{\downarrow} / \text{Speed} / \rightarrow \text{Time taken to Response}$

Autonomy → Realtime / Stock

④ → Model size → Hardware → GPU →
Electricity $\downarrow \rightarrow \text{X}$

⑤ To Shelling Log Billing ⑥



⑥ → Responsible AI / Not Biased ⇒

bias → Election → Analysis.

Hospital →

⑦ → Context Window ⇒ Stack Memory →

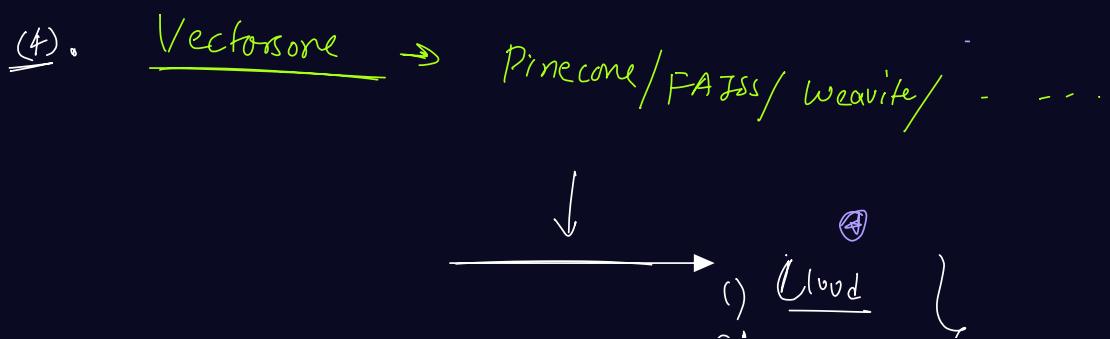
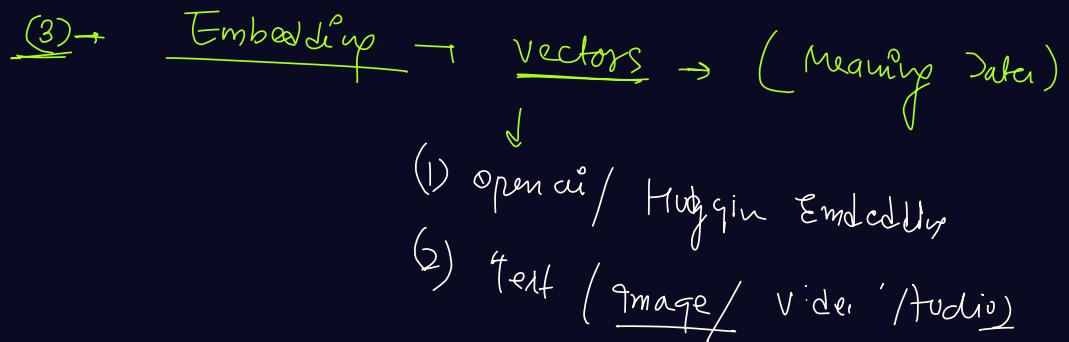
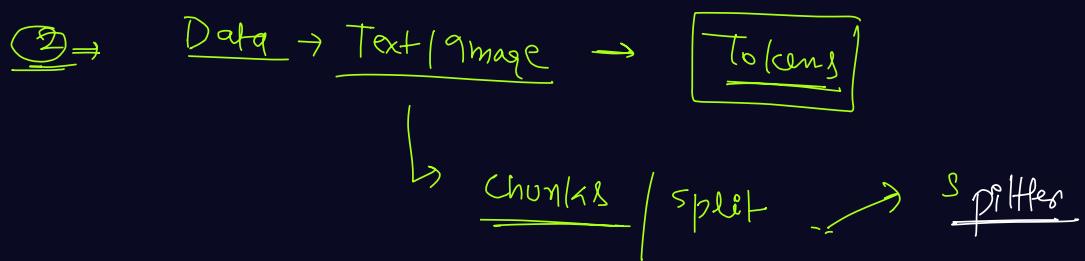
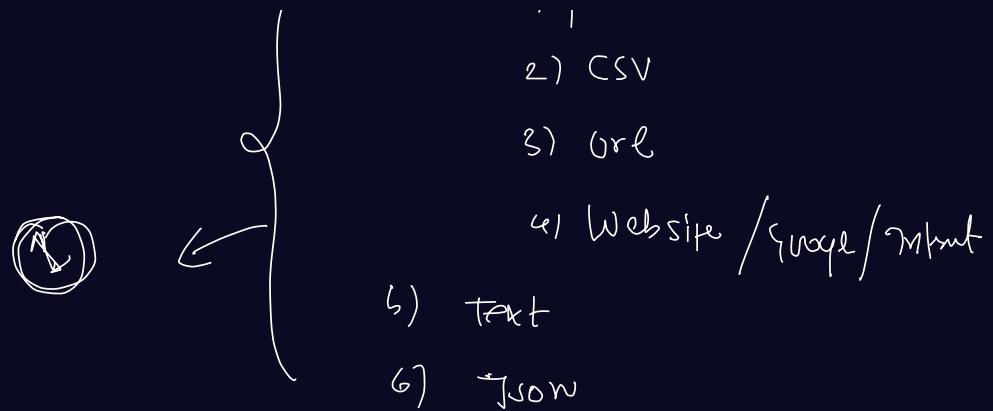
Context
↓
Memory

Memory →

⑧ Data preparation → RAG. → Langchain

(1) Data Loading →

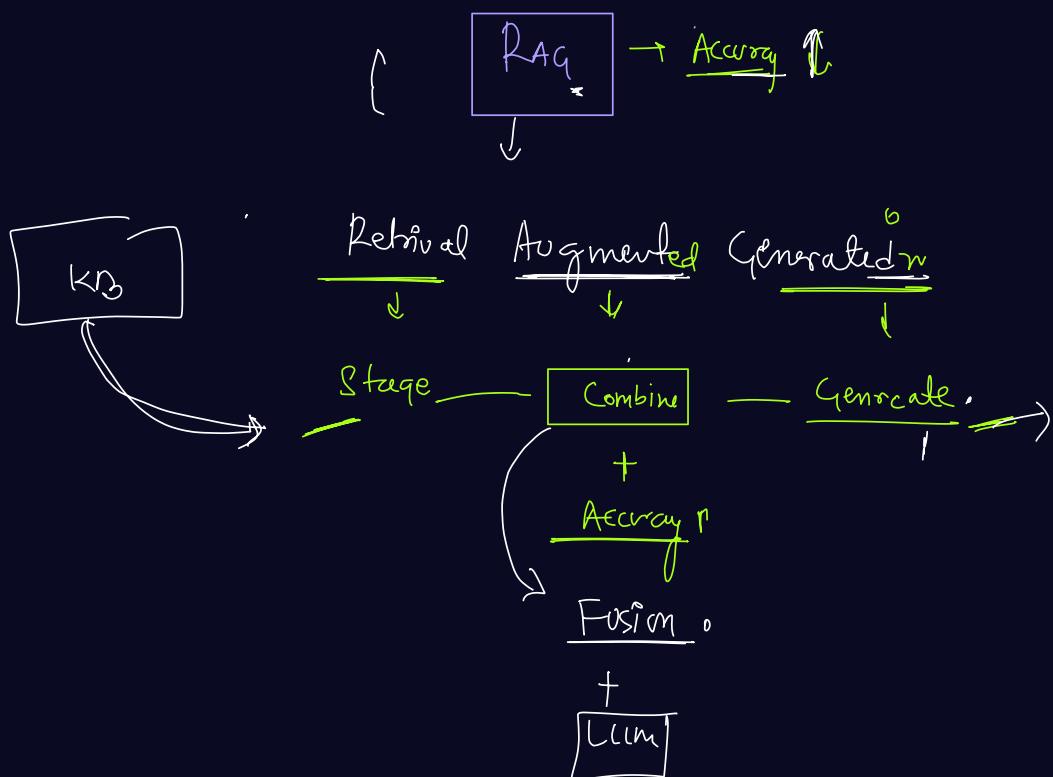
/ L) DDF



* 2) Pipeline ↗
 ④ 3) Local

↔ Data

Technique : Proprietary → prompt template }
 { RAG
 { Fine-tuning → Langchain
 } } 1) Data → Practi.





What:

- ① What is RAG
- ② Why ?
- ③ Importance ?
- ④ Use cases. ?
- ⑤ Components ?
- ⑥ Places
- ⑦ Benefit

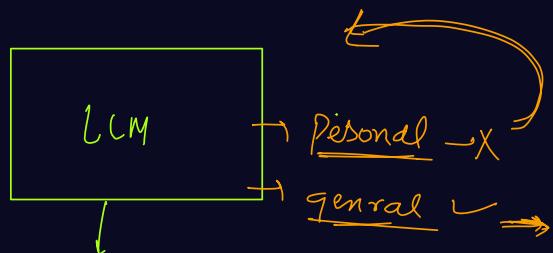
⑧) Challenges

⑨) Limitations

* Implementation

5-6

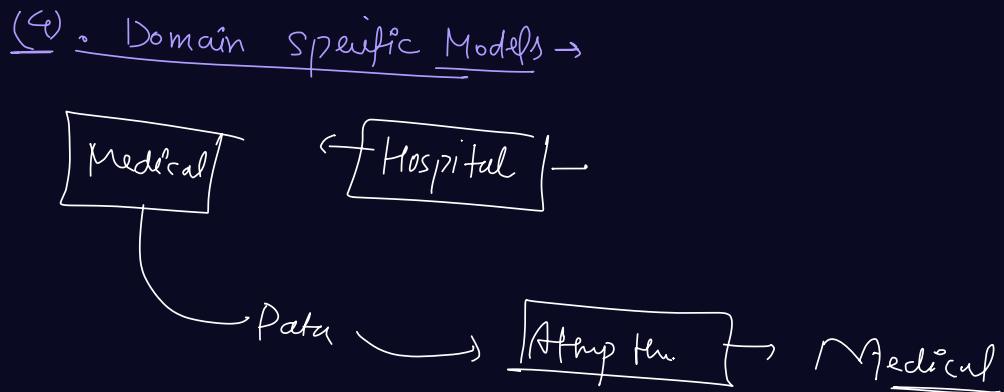
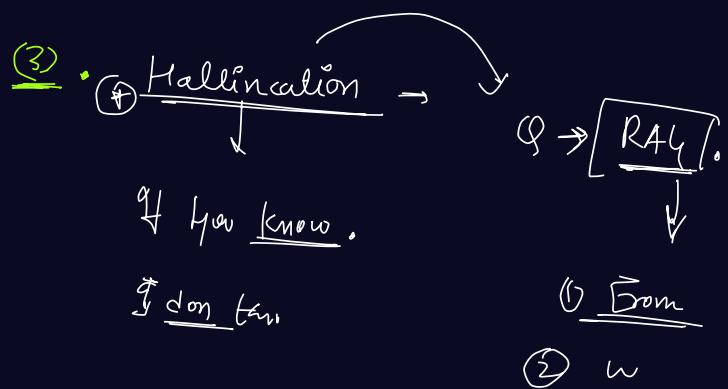
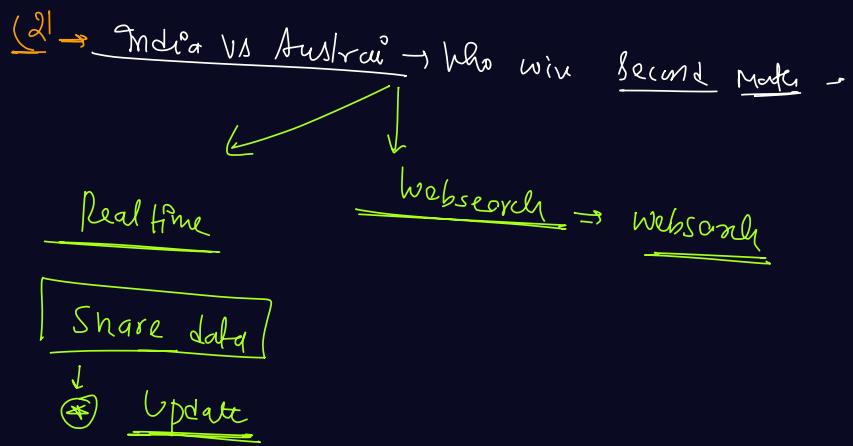
(1) → Why - ?



Q → Qnary → 9 don't like

① Personal Information

② Database



problem → general (what?) →

- - - ↓

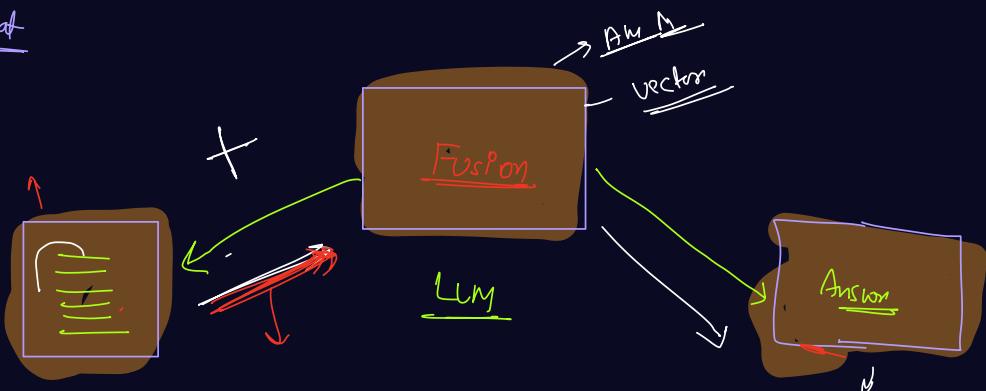
Solution

RAG

Fine tuning

QAG

What



KBS
① Personal
② Domain
③ update

Retrieval Augmented Generation

I II III

① Data → Load → chunk → embed → Vector

Agenda • 15 Dec →

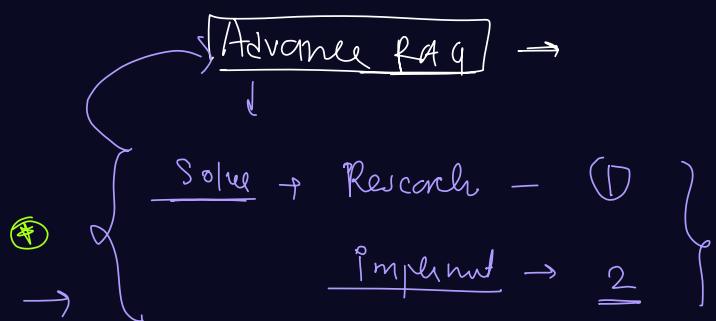
- 1) RAG
- 2) VS code setup.
- 3) ChatModel / Completion
- 4) Data Ladders

→ 5) Prompt Engineering.

⌚ Task → Langchain / OpenAI → Explode ↗

Task - 2 → 5 data ladders.

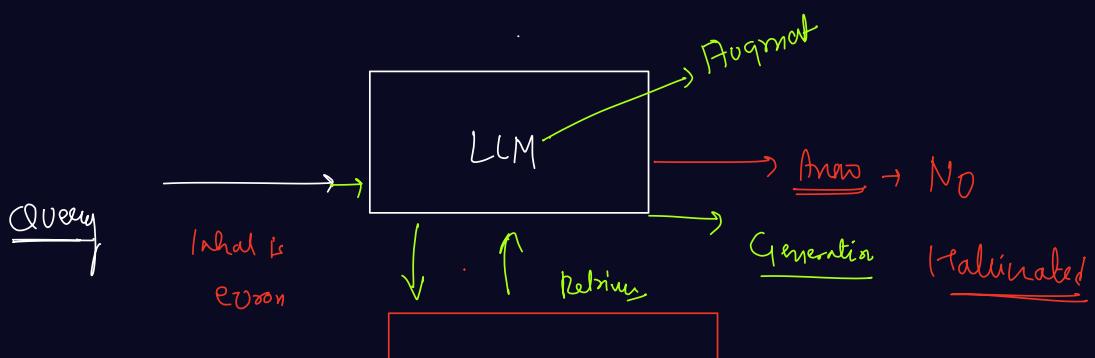
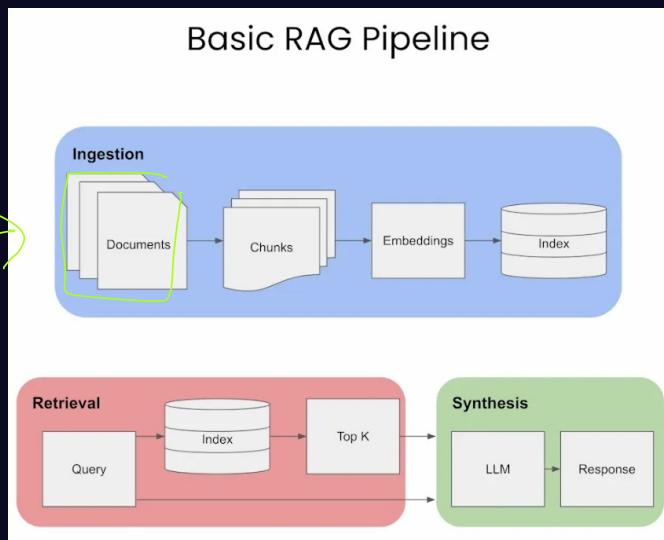
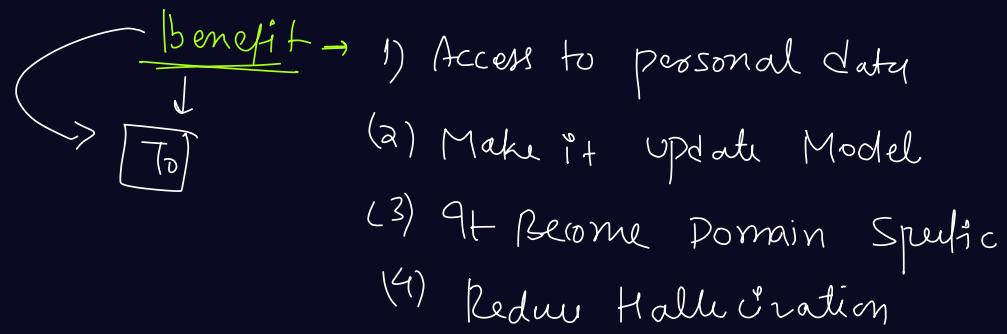
⌚ Task - 3 → RAG ↓ → Limitation ()

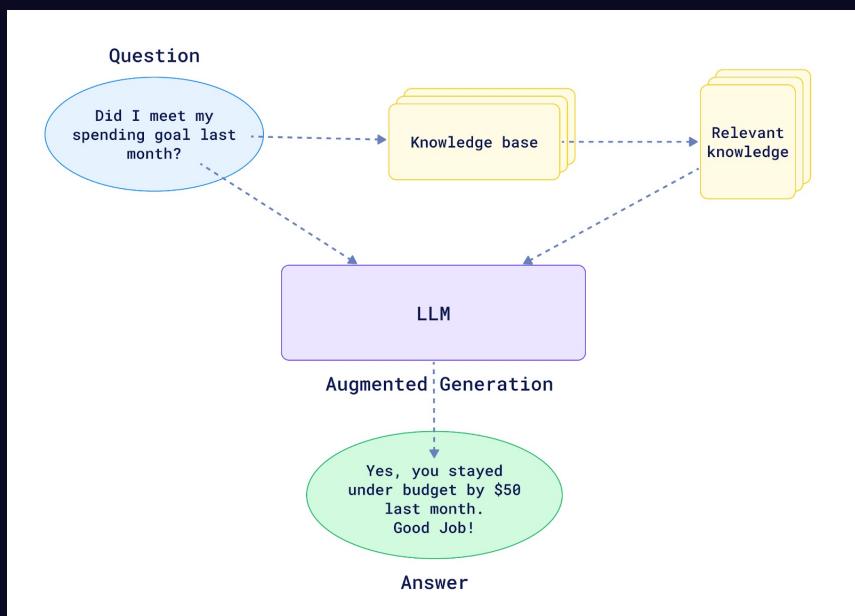
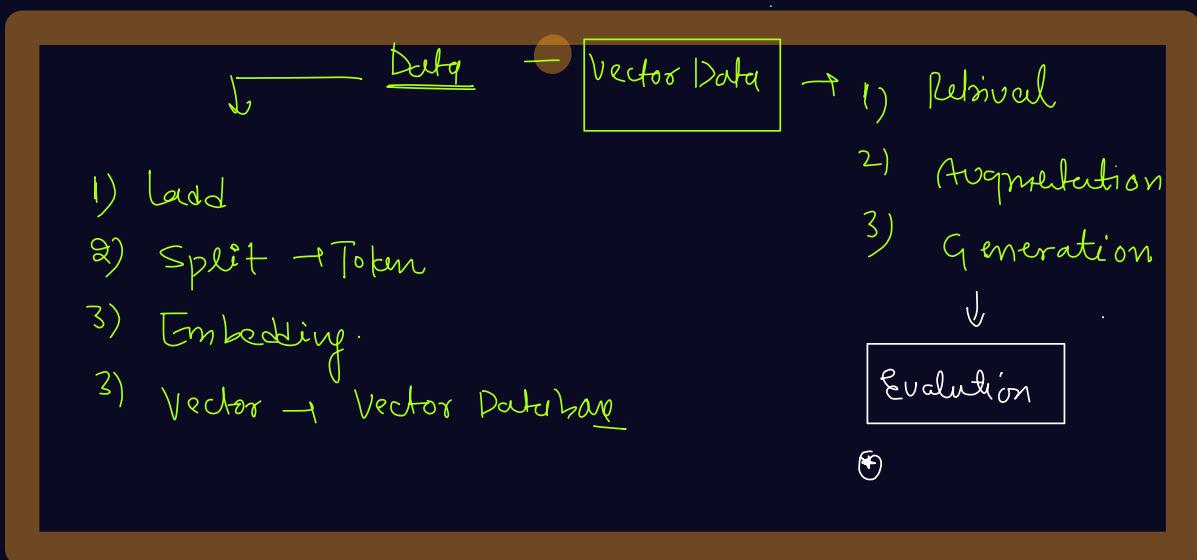
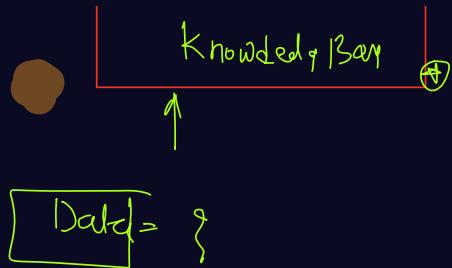


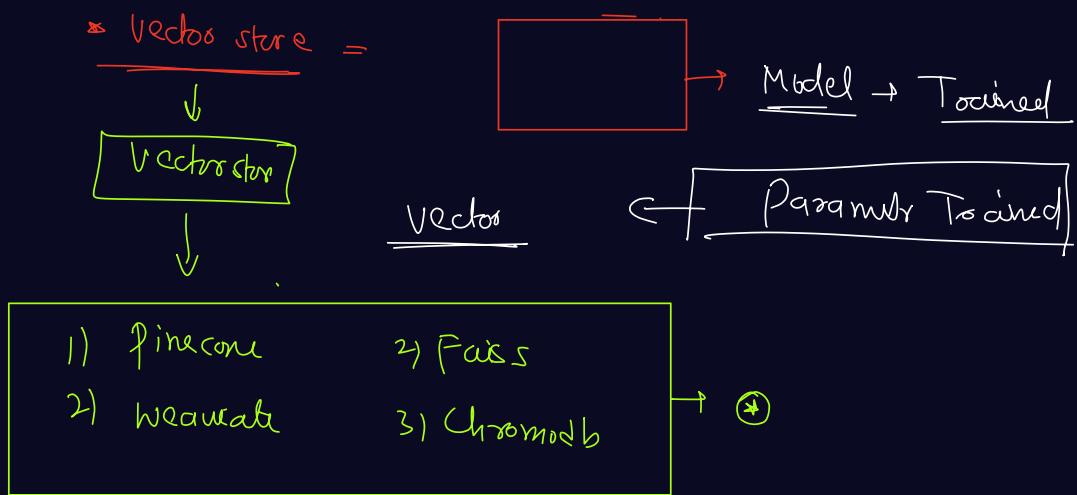
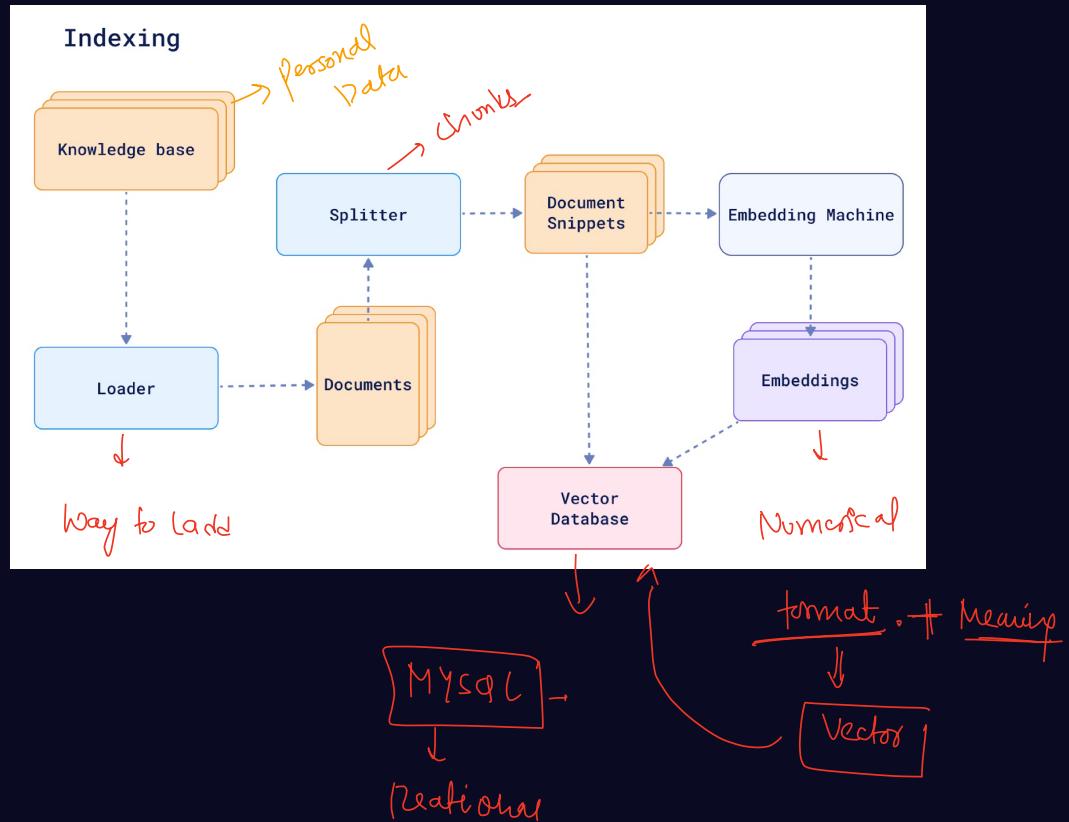
*

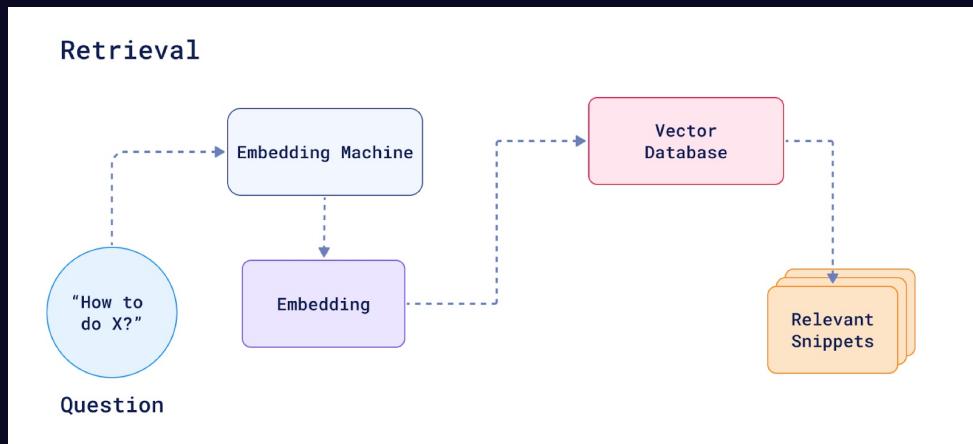
RAG → Inlay ?

In what ?

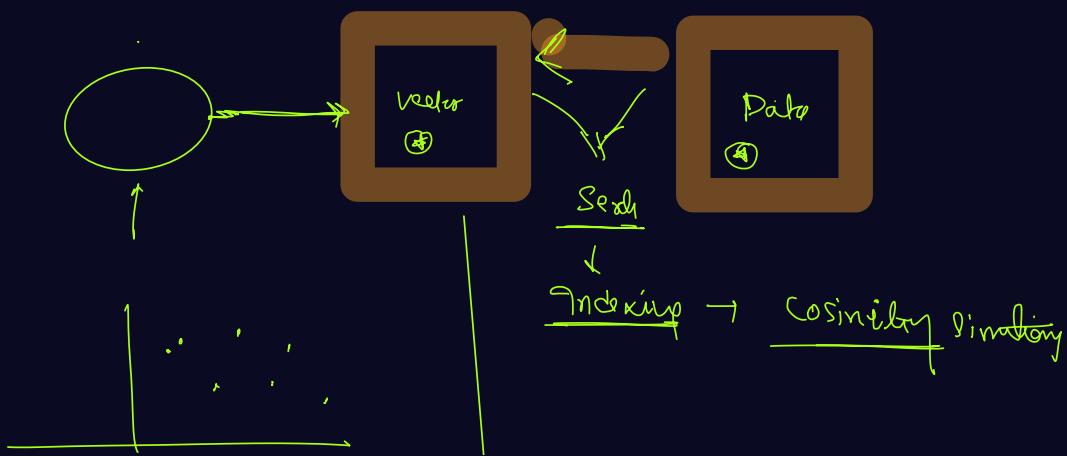
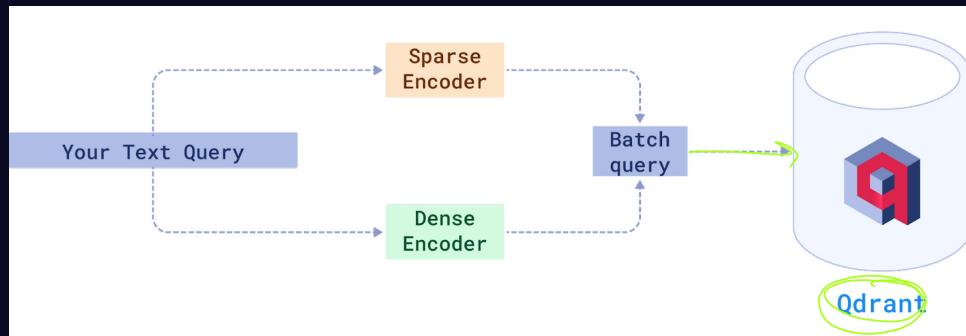


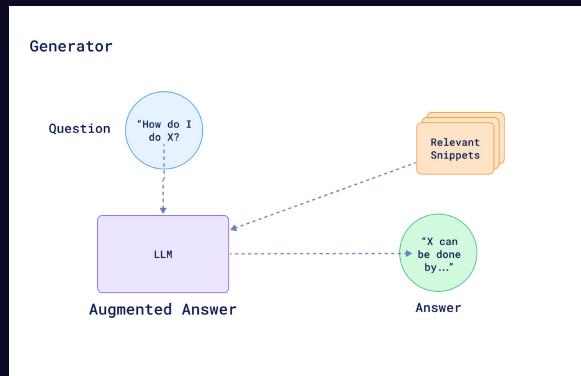
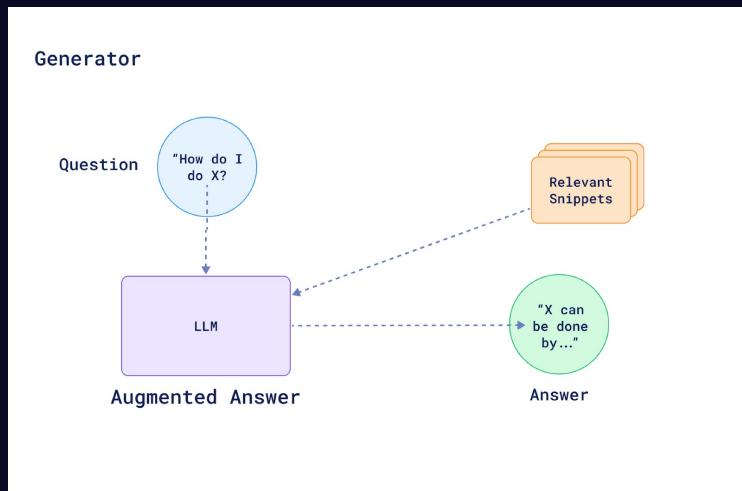
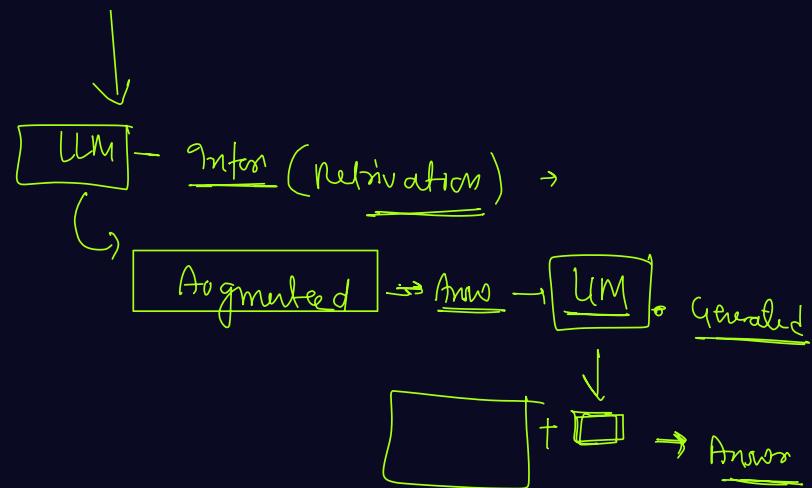


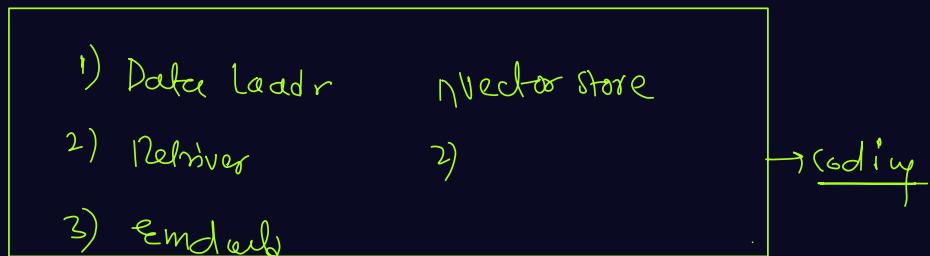
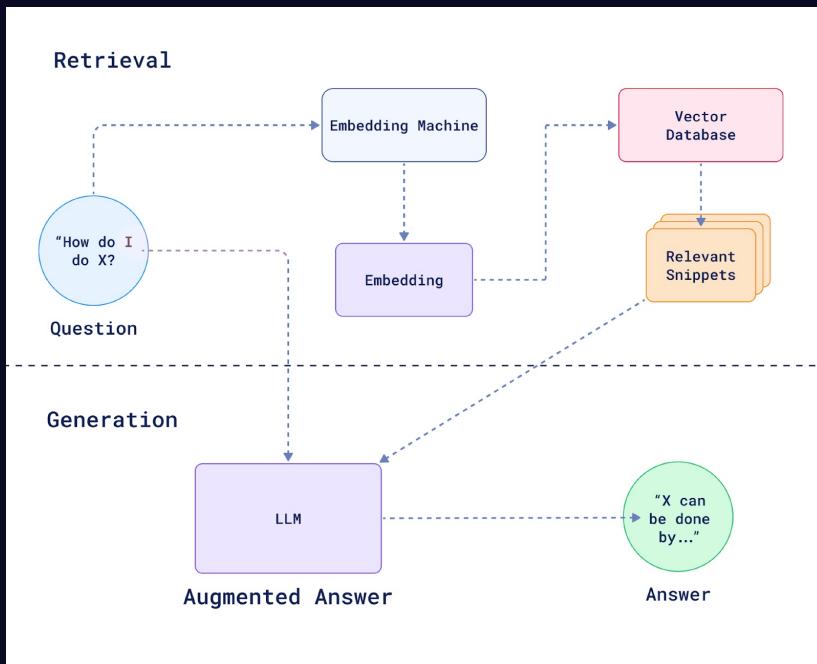




↓
What is Euro -? ① Convert Text vector







⌚ Challenges or Limitations ⇒

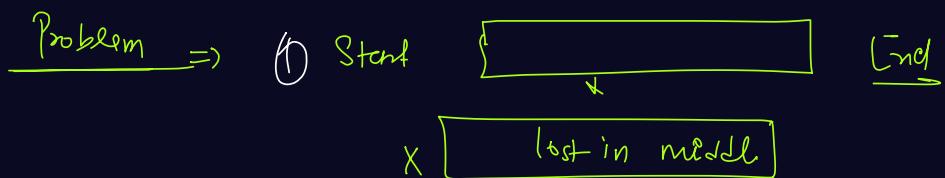
1)

- 1) Accuracy → Depend Database
- 2) Vector Datas → Storage
- 3) Realtime → update Difficult

(4).

Advanced RAG →

- ↗ {
- 1) Dense Retrieval & Hybrid Search
 - 2) Retraining . Re-ranking .
 - 3) Query Expansion
 - 4) Context Distillation



(2) Insufficient information (compression)

(3) + Complex window issue → (ST)

① Task - 1 (1) Research Advanced RAG
Technique

② Task - 2 (2) Implement RAG Architecture

(3) Task - 3 - $\sum 5$ → Modular Explor

