1) levels of Architecure → 1)
                              2)
                              3)

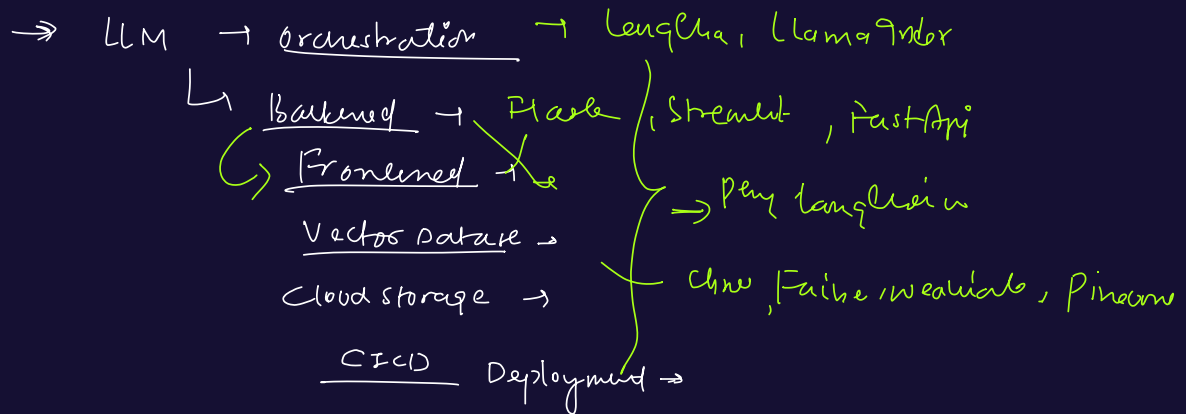(2) - Professional level LLM archичtne preview

(3) . Factors for Selecting LLM.

1) Requiomuts -
2) Architecture
3) Model size
4) Speed Latency
5) Resource & Hardware
6) Cost & ע
7) Ethical Answer

(4) - Context window          Llamq
      Multimodality            Mistral
      Accuray                  LLAMA
      Samll Psze               Blcom
                               Falcon

→ | Opensoure Model |     | Closeed Some Model |

→ LLM → Orchestration → Lengcha, Llama Index

↳ Backend → Flask, Streamlit, FastApi

↳ Frontend →

→ Peng langchain

Vector Datase →

Cloud storage →

Chue, Faihe, wealuate, Pinecone

CICD   Deployment →

---

⊕→ Open ai - pricing → Chat Model

2) Chat Compentio. Models

---

1) Level of Architecture

2) Preview of LLM Architecture

3) Tach Used in Coorse

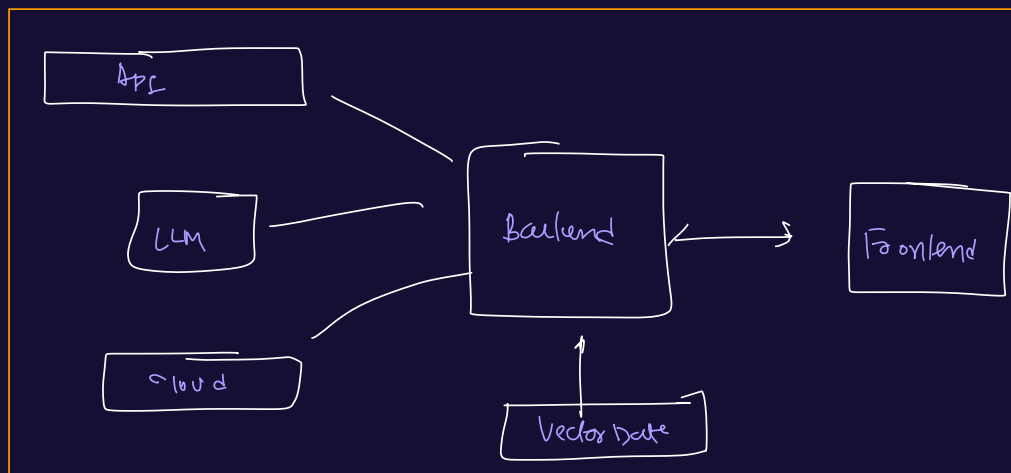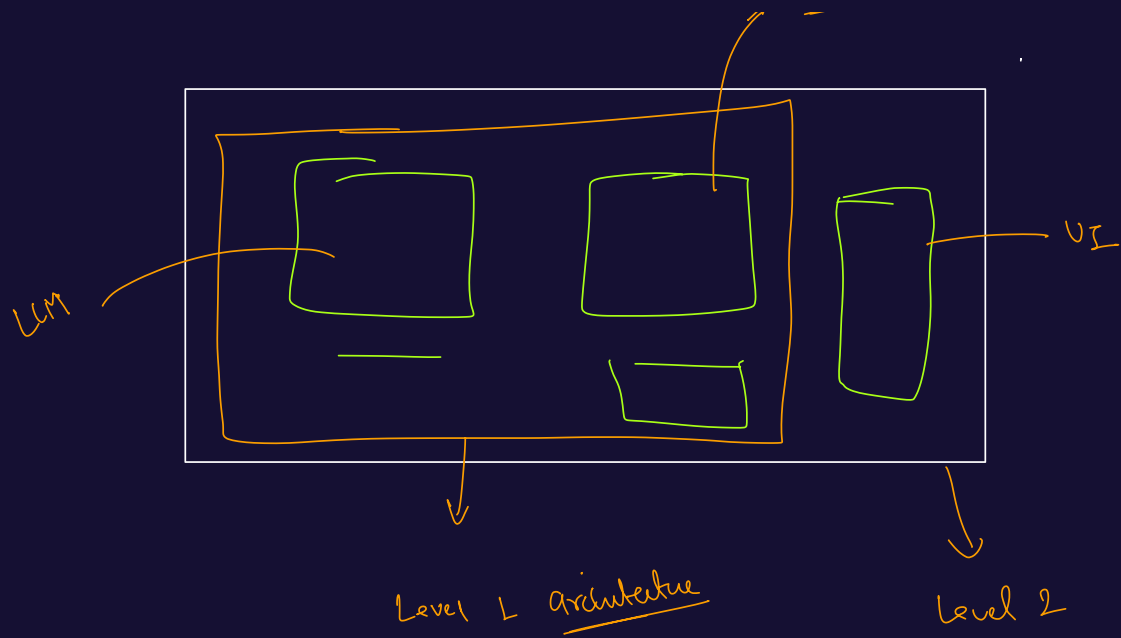4) Factors of Selecting LLMS

5) RAG Models introdtion → code.

**RAG Evaluation**

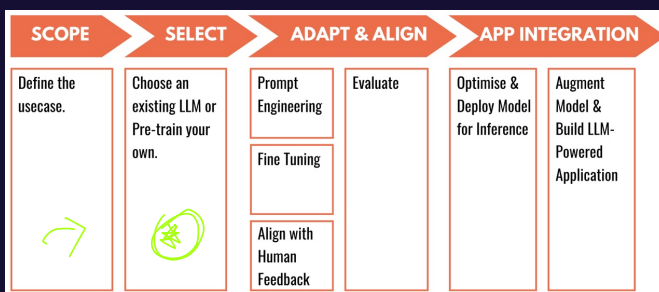|  | Generation | Retrieval |
|---|---|---|
|  | Faithfullness | Context Precision |
|  | Relevance | Context Recall |

Technique :-

⊕ Hallucination
⊛ Smaller Models
  ✳ AI Agents
  ✳ Multi Models
  ✳ Responsible
  ⊕ personalised
  ⊛ Contextual Understand

1) Models Solution

2) Fine toneing pretrained Model

3) RAG

3) Prompt Engineering

4) LLomps

Framework :

LLM

UI

Level 1 architecture

Level 2

API

LLM

Backend

Cloud

Vector Date

Frontend

1) LLM ⇒ Openai , Hugging , Llemma , GoogApI

2) Orchestration Framework → Langchain , Llama index

3) Backneal → Langchain , Langgryly .

4) Frontendend → Flask, FastApi, Strecmut

5) Vector Database → Pinecone, Chromadb , Faiss, weaviate

6) Cloud Storaqo → Gcpbucket , S3 bucket ⟷

7) Deltoyme → CICD → Github.Actions, CircleCI, Jenkensis
Docker ,

8) Multi Agents → CrewAI, Langgraph, Autogen.

ⓐ Aws Bed Rock , Vertex AI → Use LLM Globally .

| SCOPE | SELECT | ADAPT & ALIGN | | APP INTEGRATION | |
|-------|--------|---------------|---|-----------------|---|
| Define the usecase. | Choose an existing LLM or Pre-train your own. | Prompt Engineering | Evaluate | Optimise & Deploy Model for Inference | Augment Model & Build LLM-Powered Application |
| | | Fine Tuning | | | |
| | | Align with Human Feedback | | | |

1. <u>Factors</u> to <u>Select</u> → <u>LLMs</u> → → Which is <u>best</u> ?

② . <u>Future</u> trends Achievements → <u>Topic</u> Important ✪

③ . Life cycle of GenAi Model → <u>Technique</u>
   ↳

④ → <u>Tools</u> → which import — ? —① Tech stack

---

✪ <u>Future</u> → 1) Responsible AI → prompt

2) Multi Agent / Agentic AI →

3) Multimodel →

4) <u>Personalized</u> → Fine / RAG / prompt

5) <u>Contextual</u> → Memory context

---

<u>Life Cycle</u> ⇒ <u>Select Model</u>  ① → How to <u>select</u>

(2) <u>Data</u> → prompting → * <u>Load</u> / embedding / vectordatabase

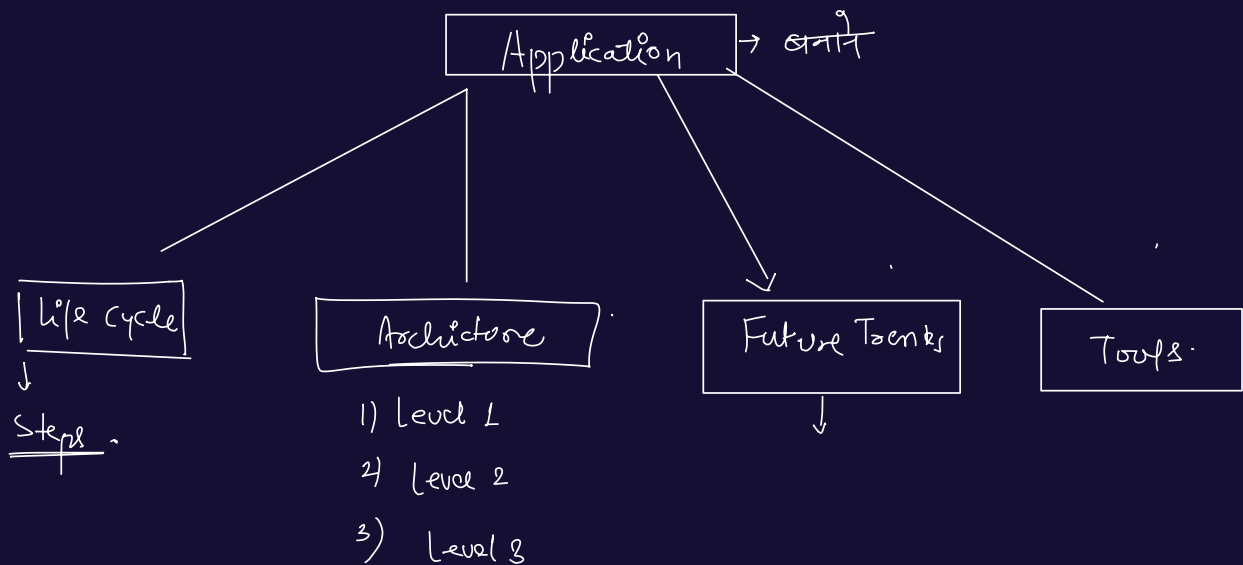⊛  (3) → <u>Technique</u> → prompt / finetune / memory / RAG .

(4) → Evaluation →

(5)

Deploy → <u>LLMops</u> → Aws Bedrock

Vertex AI

Langsmith

<u>Tech Tools</u> ⇒

Application → बनाते

Life Cycle
↓
<u>Steps</u>.

Architecture
1) Level 1
2) Level 2
3) Level 3

Future Trends
↓

Tools.

① 

Ⓣ

②

① Steps →

② Concept / Topics / Technique →

③ → Tools. / Platform.

× Sample tools. ·/ Platform .

K→ How to select LLM Models

→) Data process | Ready →

→ RAG Introduction .

How to Decide : 
→ Video Graphics
                                                                 Audio
1) Requirements. (Text, image, Multimodal)

2) Model Size

3) Accuracy → ·Hospital
   Architecture

4) Latency / Speed → Automobile

5) Hardware & Resources.

6) Cost & Licensing (openai)

7) Responsible AI - ( Not Biased) → Hospitality

8) Context Window → ·Trading

(9)

In the openai

① $ Ecosystem
   |
   +
② Connection with
   Others

★ Models → open Vs Closed Source
                 ↓              ↓
        ↳ Free           Closed

Response :  ← 1) Laama      1) Gpt 4
   ↓
Error        2) Mistral     2) Claude

                            3) Gemini

3) Falcom

4) Cohere
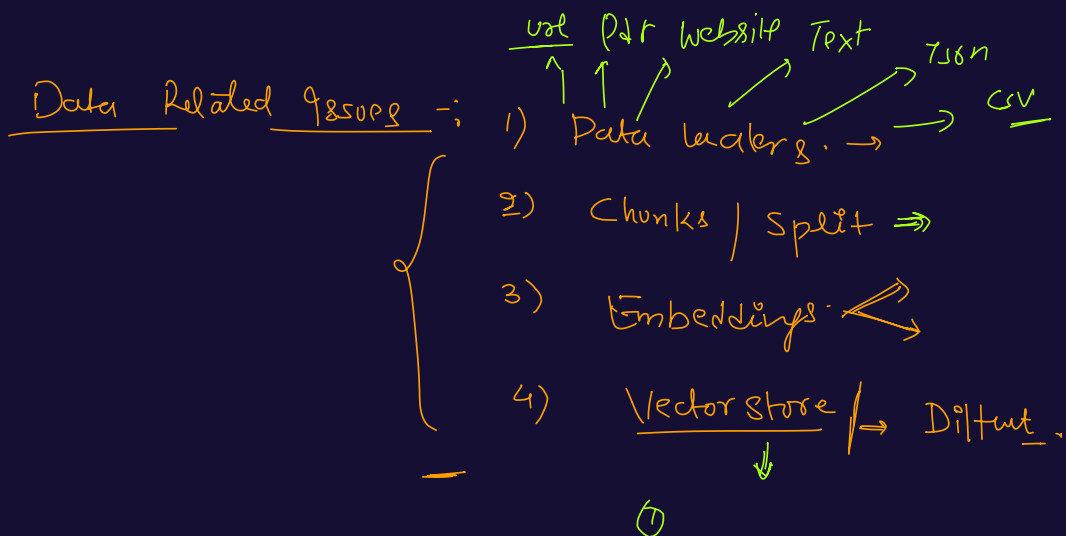
4) Bloom .

③ 🌟 Accuracy , Latency, Ecosystem ⇒

④ -How to use → ① API

② Download

③ Cloud Based → Aws Bed Rock / GCP

Vertex AI

Data Related Issues -:

```
                    url Pdr website Text    Json
1) Data leaders. →              → csv
2) Chunks / Split ⇒
3) Embeddings. ↙
4) Vector store /→ Ditturt.
           ⇓
           ①
```

| ① * RAG | ② Finetuning | ③ * Prompting |

① At Agent

② → — RAG.

1)