

# RAG

1 What is diff b/w Simple AI Assistant vs RAG

2 Finetuning & RAG

3 Architecture of RAG

4 Implementation of RAG

5 Langchain -

6 Vector Database -

7 Multimodal RAG

8 Graph RAG

9 CAG → fast retrieval

↑  
Cache

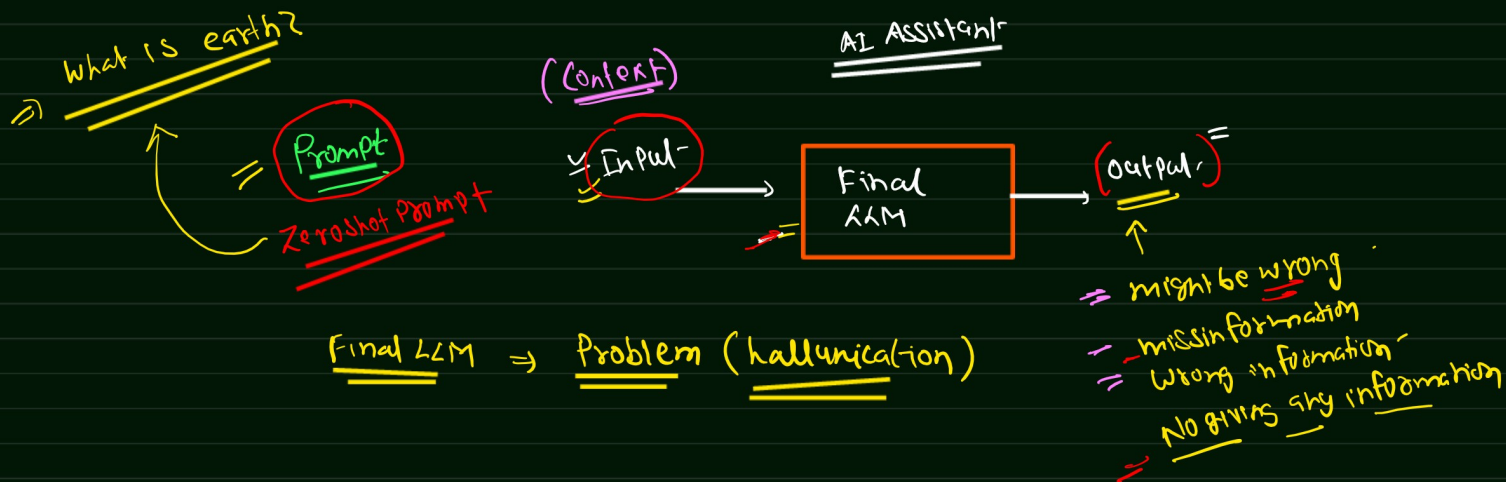
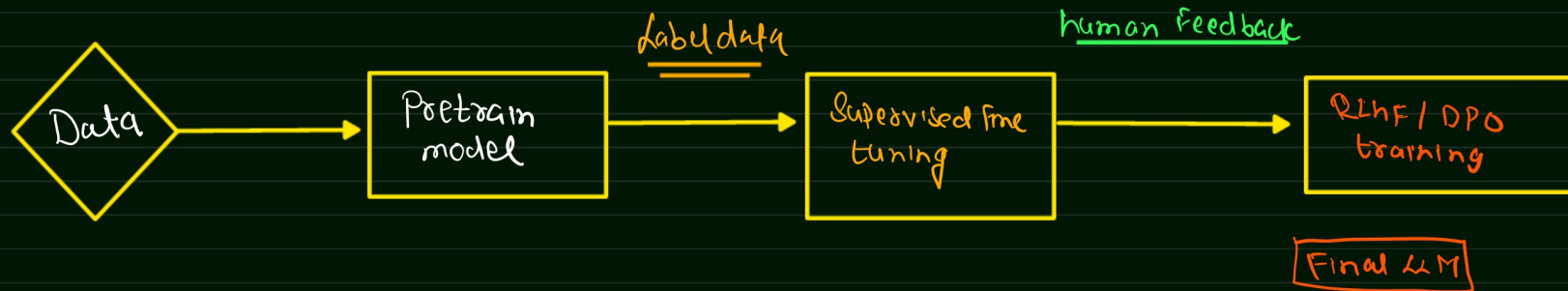
- 10 Agentic RAG =

RAG ⇒ look & sample  
↳ variable

# Diff B/w AI Assistant and RAG

## Retrieval Augment- generation

- ↳ ① Data ingestion.
- ↳ ② Data Retriever
- ↳ ③ Data generation



Prompt ⇒

- 1 ZeroShot Prompt
- 2 Few Shot Prompt
- 3 Chain of thought Prompting

Agent ← different instructions } prompting  
Input

Internet data → LLM  
↳ Entire data from data

(Medical domain) → XYZ ⇒ Researching on one specific medicine (latest) Cancer

Jaspreet ⇒ Company ⇒ (marketing sales)

Q12

1.f

LLM

retrain

RA4

Context Few-shot learning

Finetuning  
SFT, DPO

⇒ Few-shot Prompting ← RA4

{ Arctical }  
(refer above arctical and tell me what is earth?)

# Simple AI Assistant

Input  
(Prompt)

LLM

Output

hallucinating ⇒ misleading info  
wrong output

Current info  
or any info which is  
not present on internet  
any company specific info  
any domain specific info

expensive

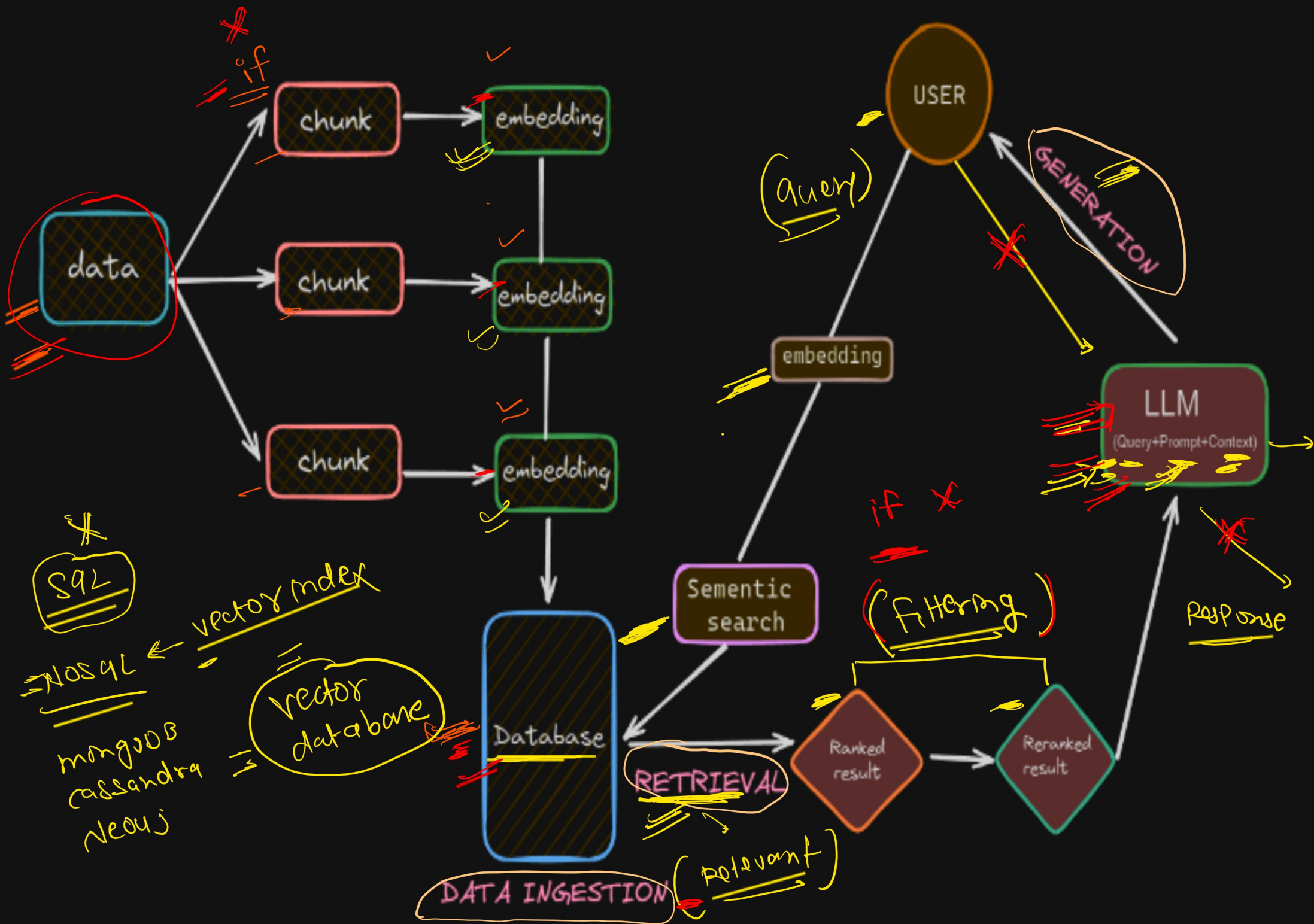
① First solution is :- Finetuning (Retraining)

STF, DPO

Research

very very popular

② RAG





✓ { OPENAI API ✓  
Gemini API ✓  
GROQ API ✓  
Claude API ✓  
Huggingface API ✓ } key

Langchain

~~Langchain~~ Langchain

Langchain-Groq

Langchain-Gemini

Dev