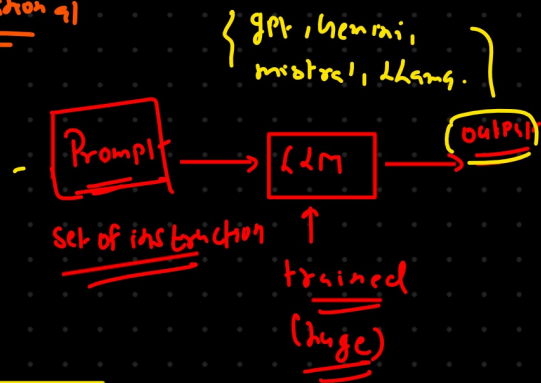


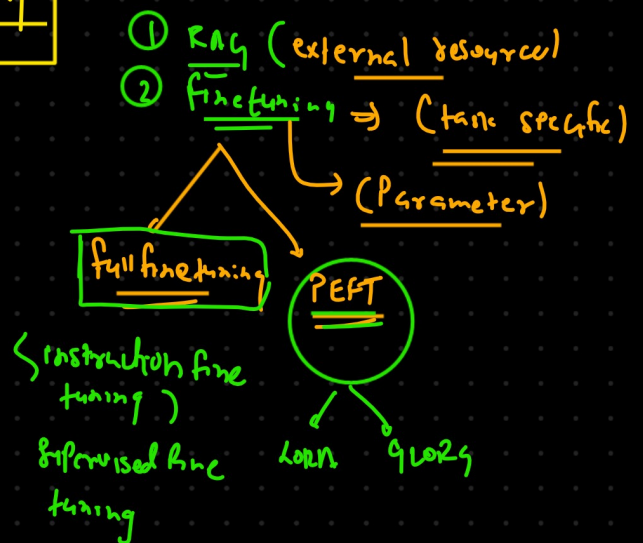
- fundamental
- ① Roadmap
 - ② Intro
 - ③ Google Gemini, hugging face, openAI API
 - Google AI studio
 - token
 - GPT 3.5, GPT 4

NLP Related Foundation

- ① text cleaning
- ② encoding & embedding
- ③ transformer
- ④ BERT API



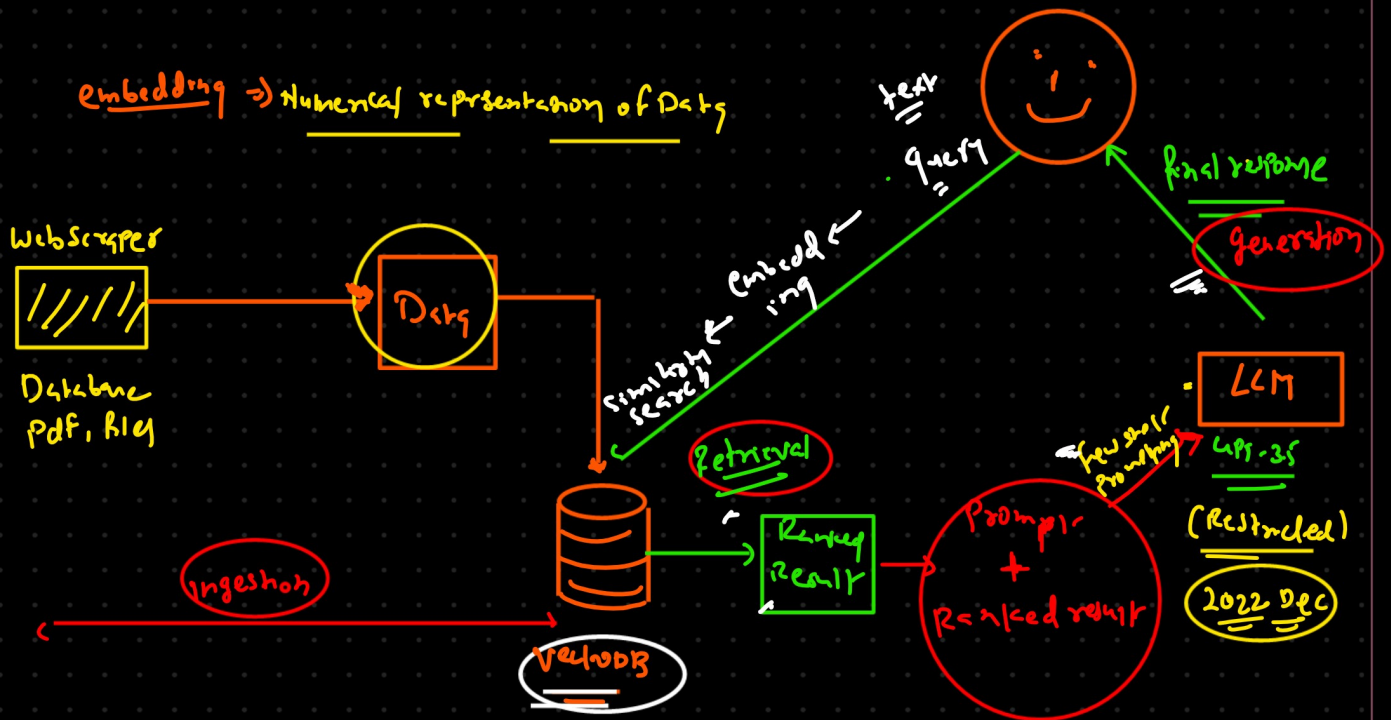
text cleaning



- ③ RLHF
- ④ fact checking

Cosine Similarity Dot Product Jaccard Similarity
Embedding \Rightarrow $\begin{bmatrix} \text{ } \end{bmatrix}$ $\begin{bmatrix} \text{ } \end{bmatrix}$ $\begin{bmatrix} \text{ } \end{bmatrix}$ Proof

embedding \Rightarrow Numerical representation of Data



RAG

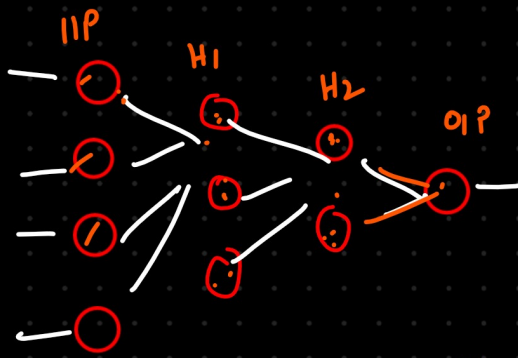
- 1 Ingestion
- 2 retrieval
- 3 generation

RAG \Rightarrow Retrieval Augmented Generation

RAG \Rightarrow 8110

- VectorDB
- similarity
- length or 26 index
- Retrieval
- F
- Combined
- hybrid

fine-tuning \Rightarrow technically



$$4 \times 3 = 12$$

$$3 \times 2 = 6$$

$$2 \times 1 = 2$$

$$12 + 6 + 2 = 20$$

$$15 + 8 + 3 = 26$$

$$\Rightarrow 26$$

Parameters = 26

↑
trainable
param

Weights + bias

Basic Unit

- ↳ CNN
- ↳ LSTM
- ↳ Transformer
- ↳ GNN
- ↳ RL

- ① FF (Dot Prod + Act + Loss)
- ② BP (optimizer)

↳ LD, Adam, RMS Prop, Adagrad

minimizing Loss ↓ ↓ ↓

100 billion

→ 1.76 Trillion (Size of model + Data)

Model	Provider	Open-Source	Speed	Quality	Params	Fine-Tuneability
GPT-4	Open AI	<input checked="" type="checkbox"/>	☆☆☆	★★★★	1.76T	<input checked="" type="checkbox"/>
GPT-3.5 Turbo	Open AI	<input checked="" type="checkbox"/>	★★★★	★★★★	175B	<input checked="" type="checkbox"/>
GPT-3	Open AI	<input checked="" type="checkbox"/>	☆☆☆	★★★★	175B	<input checked="" type="checkbox"/>
Claude 3 Opus	Anthropic	<input checked="" type="checkbox"/>	☆☆☆	★★★★	Not specified	<input checked="" type="checkbox"/>
Claude 3 Sonnet	Anthropic	<input checked="" type="checkbox"/>	☆☆☆	★★★★	Not specified	<input checked="" type="checkbox"/>
Claude 3 Haiku	Anthropic	<input checked="" type="checkbox"/>	☆☆☆	★★★★	Not specified	<input checked="" type="checkbox"/>
Command Nightly	Cohere	<input checked="" type="checkbox"/>	☆☆☆	★★★★	52B	<input checked="" type="checkbox"/>
BERT	Google	<input checked="" type="checkbox"/>	★★★★	☆☆☆☆	345M	<input checked="" type="checkbox"/>
T5	Google	<input checked="" type="checkbox"/>	☆☆☆	☆☆☆☆	11B	<input checked="" type="checkbox"/>
PaLM	Google	<input checked="" type="checkbox"/>	☆☆☆	☆☆☆☆	540B	<input checked="" type="checkbox"/>
Gemini Nano	Google	<input checked="" type="checkbox"/>	☆☆☆	☆☆☆☆	1.76T	Coming soon
Gemini Pro	Google	<input checked="" type="checkbox"/>	☆☆☆	☆☆☆☆	175B	Coming soon
Gemini Ultra	Google	<input checked="" type="checkbox"/>	☆☆☆	★★★★	175B	Coming soon
LLaMA	Meta AI	<input checked="" type="checkbox"/>	☆☆☆	☆☆☆☆	65B	<input checked="" type="checkbox"/>
Llama 2 7B	Meta AI	<input checked="" type="checkbox"/>	☆☆☆	☆☆☆☆	7B	<input checked="" type="checkbox"/>
Llama 2 13B	Meta AI	<input checked="" type="checkbox"/>	☆☆☆	☆☆☆☆	13B	<input checked="" type="checkbox"/>
Llama 2 70B	Meta AI	<input checked="" type="checkbox"/>	☆☆☆	★★★★	70B	<input checked="" type="checkbox"/>
Mistral 7B	Mistral AI	<input checked="" type="checkbox"/>	☆☆☆	★★★★	7.3B	<input checked="" type="checkbox"/>
Mixtral 8x7B	Mistral AI	<input checked="" type="checkbox"/>	☆☆☆	★★★★	46.7B	<input checked="" type="checkbox"/>
Orca	Microsoft	<input checked="" type="checkbox"/>	☆☆☆	★★★★	13B	?
Falcon 40B	Falcon LLM	<input checked="" type="checkbox"/>	☆☆☆	★★★★	40B	<input checked="" type="checkbox"/>
Falcon 180B	Falcon LLM	<input checked="" type="checkbox"/>	☆☆☆	★★★★	180B	<input checked="" type="checkbox"/>

$NN \rightarrow \text{Weight} + \text{Bias} \Rightarrow \text{Parameter}$
 \uparrow (BP)
trainable Parameter

LLM

$\left\{ \begin{array}{l} \text{GPT} \Rightarrow 175B \\ \text{Llama} \Rightarrow 70B \\ \text{mistral} \Rightarrow 45B \end{array} \right\}$
 ;

Fine-tuning \Rightarrow Retrain the Parameter of Network (model)

Fullfinetuning
 Set of Parameter
(PEFT)

(Q)
RAG -
 external Data source \Rightarrow eng.
 (45/100)
 (Donor)
Fine-tuning
Weight + bias
(Resource)
 (10/100)

Combined approach

AI

Fine-tuning \Rightarrow RAG APP

