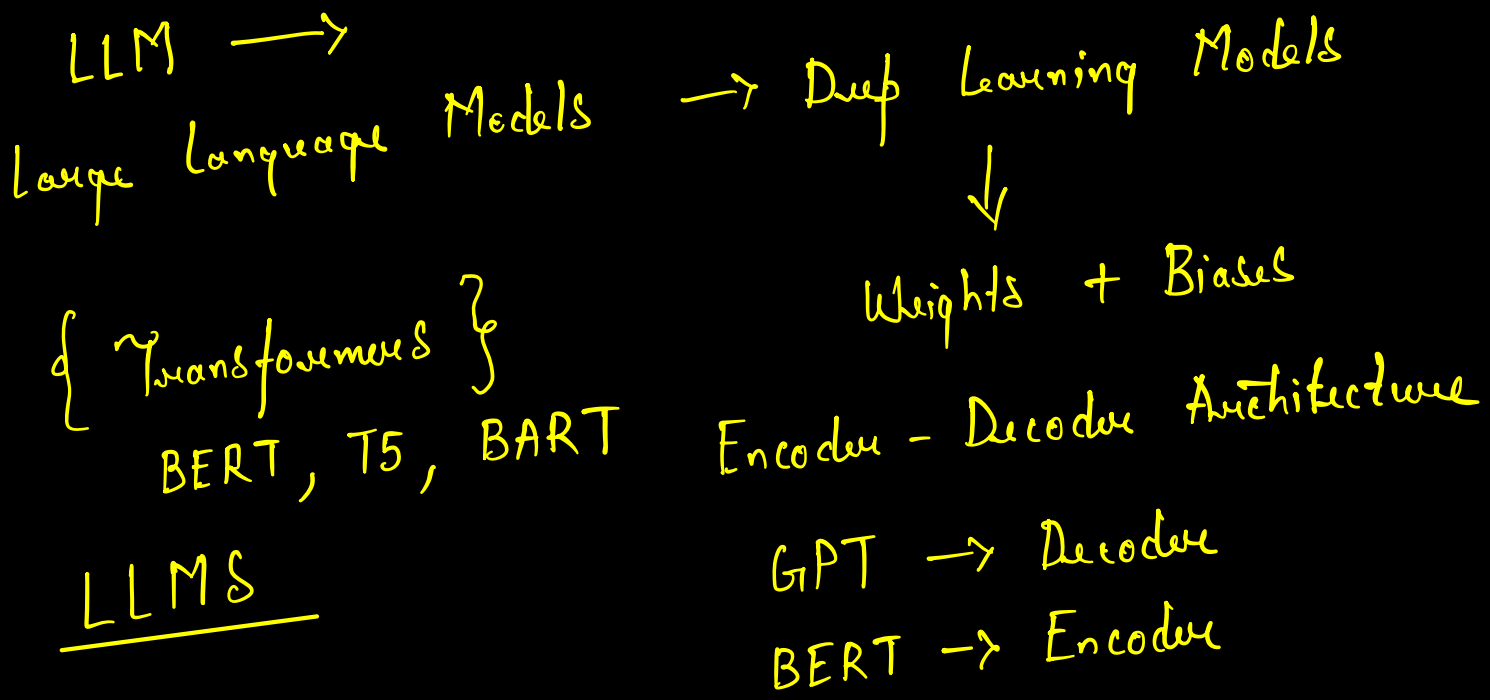


Introduction to RAG



NLP

Upstream :- Pretraining (Masking)

Downstream :- Post Training

↓
Fine Tuning

Text Classification, Summarization, POS, NER
QA, MT, Text Generation

LLM Problems

Base Model, Foundational Model

SLM

Small Language Model

+
{ Custom Data }



Fine Tuning Task

1) Knowledge Cutoffs

↳ Cutoff Training Date

2) Lack of Info
in Domain specific Task

Wrong Answer

3) Lack of Private Data

4) Loses Trusted Source

5) Probabilistic Output Generation.

{ RAG }

vs

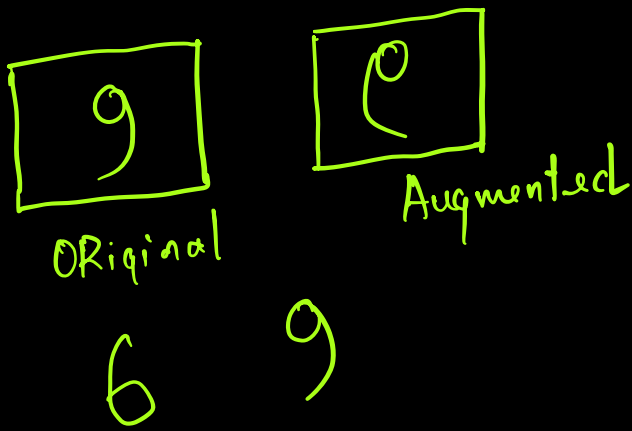
{ LLM }

✓✓

Training { limited }
{ Computationally Expensive }
GPU

Basics of RAG

Retrieval Augmented Generation → LLM
(Original → Modification)


Original Augmented
6 9

My name is Paul.
Paul is my name.
99% Context
+
Modification

Data → Vector DB → * Sources Original
↓
LLM
↓
* Answer

{ Context Retained }

RAG Simple Architecture

1) Data Ingestion (External Knowledge Sources)

Text \rightarrow PDF, HTML, JSON, CSV, XLX
DOCX, PPTX

MultiModal $\left\{ \begin{array}{l} \text{Image} \rightarrow \text{GIFS, JPG, PNG, TIF} \\ \text{Audio} \rightarrow \text{.stt, .mp3, .acc} \\ \text{Videos} \rightarrow \text{.mp4, .mov} \end{array} \right.$

2) Vector Databases

Number Representation

My name is Paul.
0.1 0.8 0.9 0.75

$[0.1, 0.8, 0.9, 0.75] \rightarrow$

Algorithm :- Embeddings

1) Bag of Words

- 2) TF IDF
- 3) Cooccurency Matrix
- 4) Word 2 Vec $\begin{cases} \text{CBOW} \\ \text{Skipgram} \end{cases}$
- 5) (Positional Encoding) *
- 6) One Hot Encodding
- 7) FLMO

1) Text Data / Multimodal

↓
Numbers (Embeddings)

(Model)

→ MTEB

1) Open Source → Hugging face

gemini-embedding model
all mini lm 6

2) Paid →

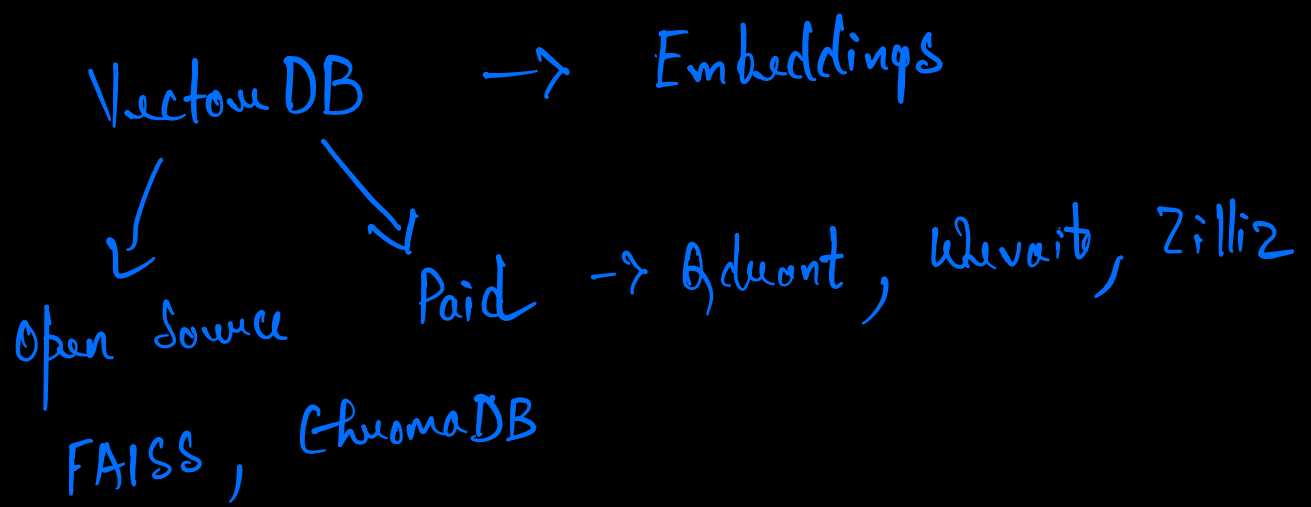
OpenAI

Ti Tan

Algorithm :- Text / Image

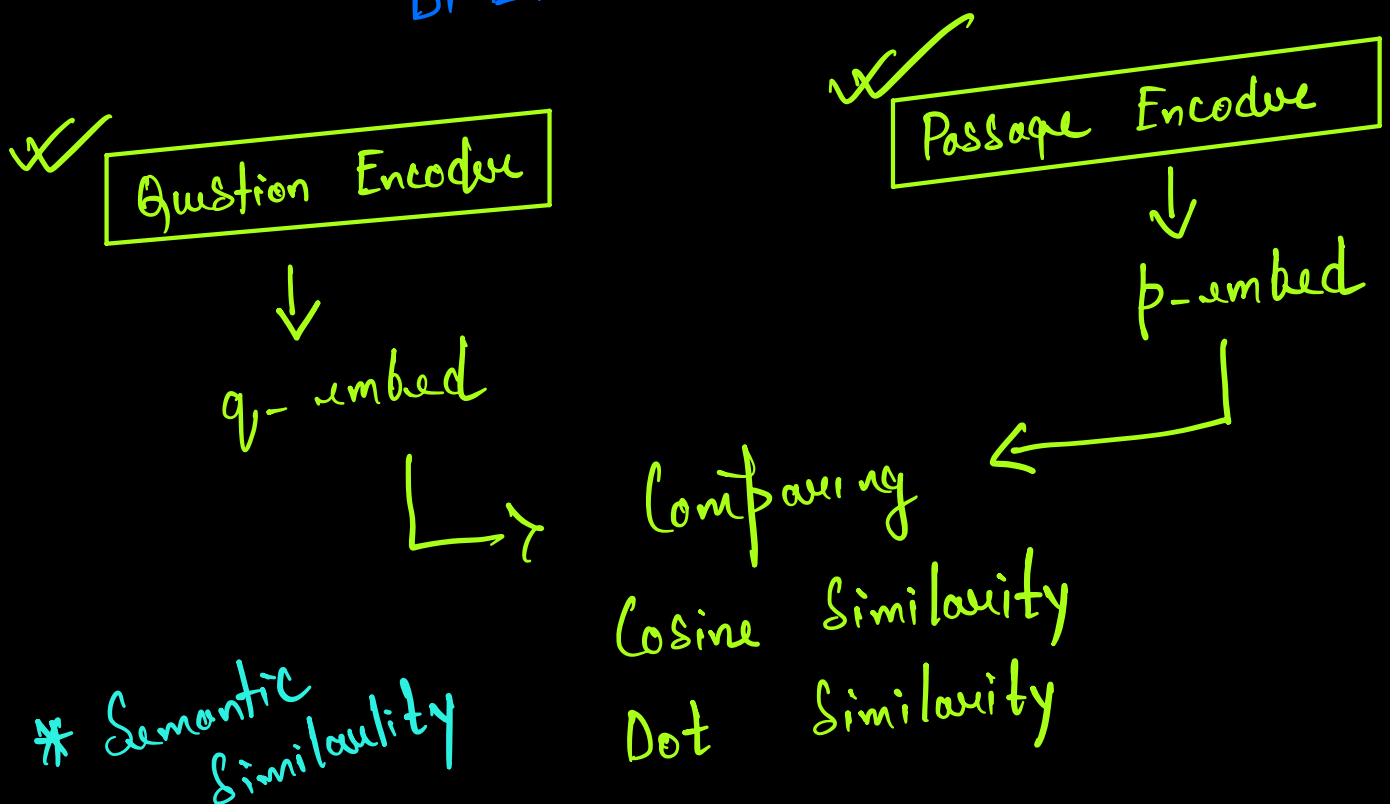
↓

{ Numeric Representation }

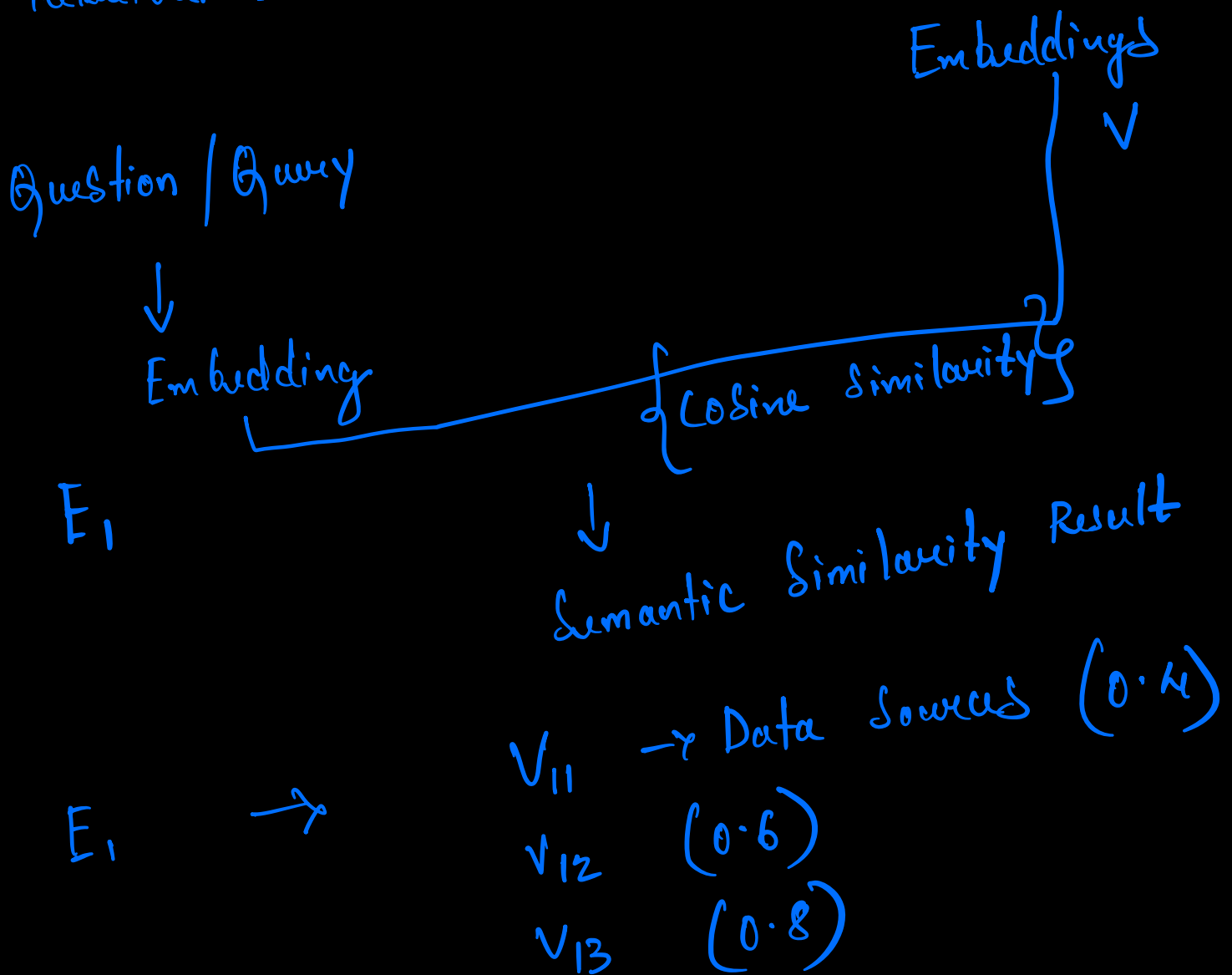


3) Retrieval
Neural Retrieval DPR → 2 Encoders
Dense Passage Retrieval

Bi Encoders



Retrieval :- Fetch relevant data from VectorDB



4) Augmentation

- 1) User Query
- 2) Retrieved Docs
- 3) Prompt

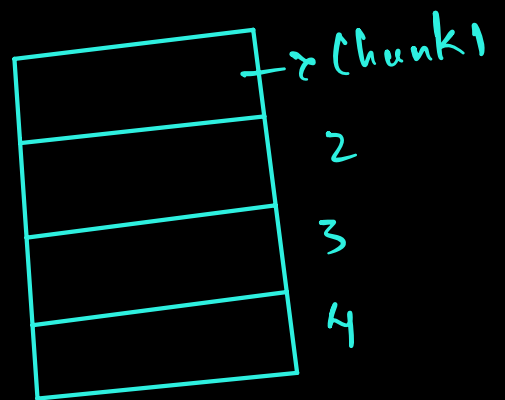
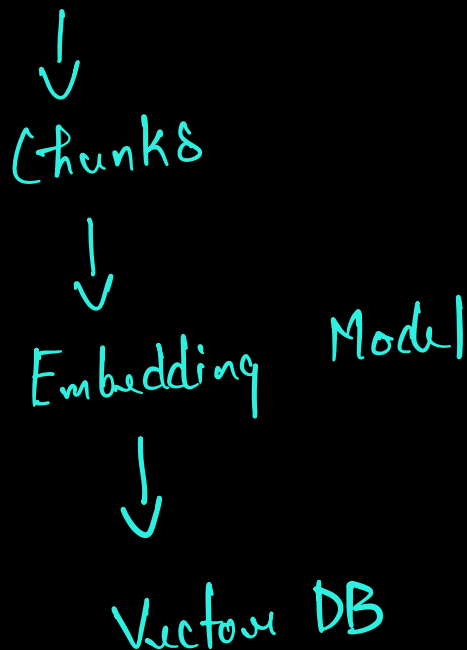
5) Generation

LLM \rightarrow Text Generation

\rightarrow Context (3 parts)
Q, R, P

Phase 1 :- Data Ingestion

1) Docs, Crawl



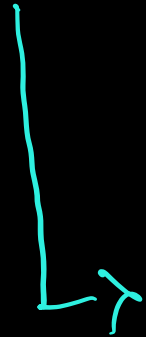
Phase 2 :- Retrieval

Query → Embedding Model → Vector DB

{ Dense Vectors }
{ Sparse Vectors }

Retrieved
Docs

* Reranking



Context + Query + Prompt

↓
LLM

↓
Response

↓
Set of
instructions

(

Introduction

Application

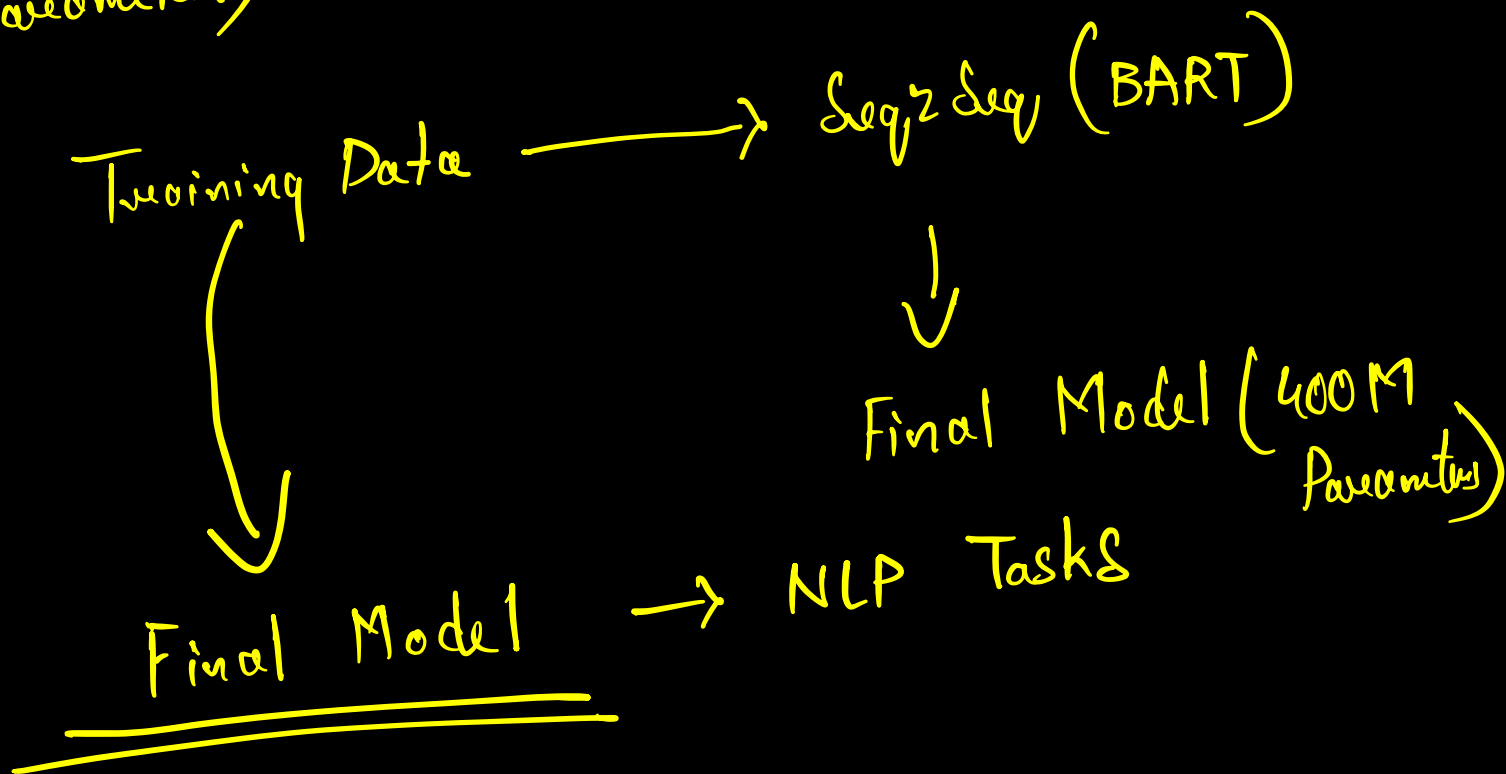
Use Case

Real World

Conclusion

- 1) Extractive :- Exact same piece of info from data
- 2) Abstractive :- Modified info but with same context. like Rephrasing.

Parametric Pretrained Seq2Seq (BART)
(400M Parameters) that stores knowledge in weights



Parametric Memory

Non Parametric Memory

Dense Vectors of any type of external
Index data sources
(Wikipedia)

Accessed : - Neural Retrieval
DPR

2 RAG Recipes

- 1) RAG Sequence :- Uses the same retrieved docs to generate the entire output sequence
- 2) RAG Token :- We can use different documents to generate the final output sequence.