

Retrieval Augmented Generation (RAG)



Large language Models

Transformer based Neural Networks

LLMs: huge number of parameters

weights & biases

parametric knowledge



Wikipedia, GitHub, Internet

Challenges (Vanilla LLM):

pretrained on a large Corpus of textual data

1) Prompt "How old is Harry?"

Harry Potter

"What are procedures that are followed in Amazon?"

2) knowledge cut off date

Pretraining LLM:

GPT-4 GPT-3

~~BB~~

Oct 2024

2025

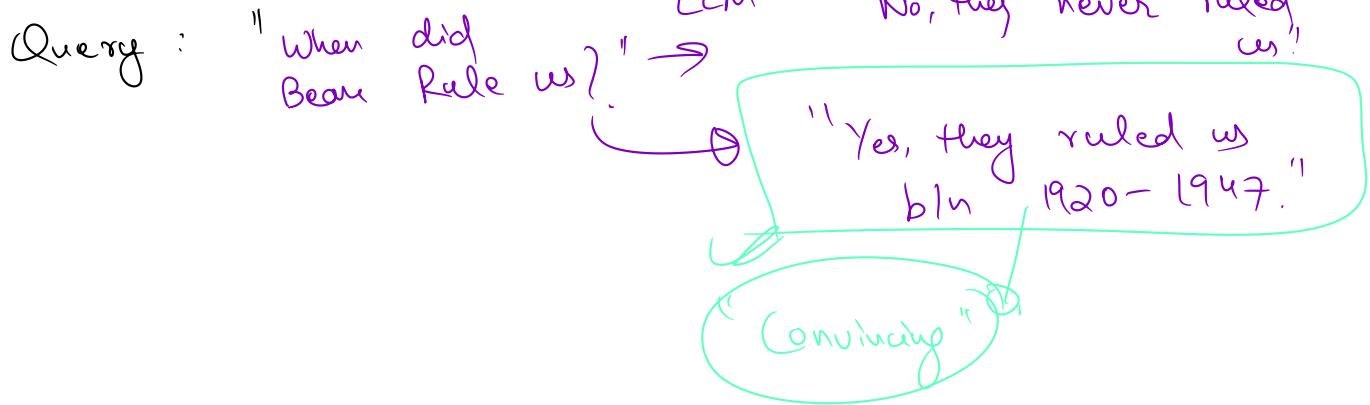
13th Jan 2026

LLM has no access to the data

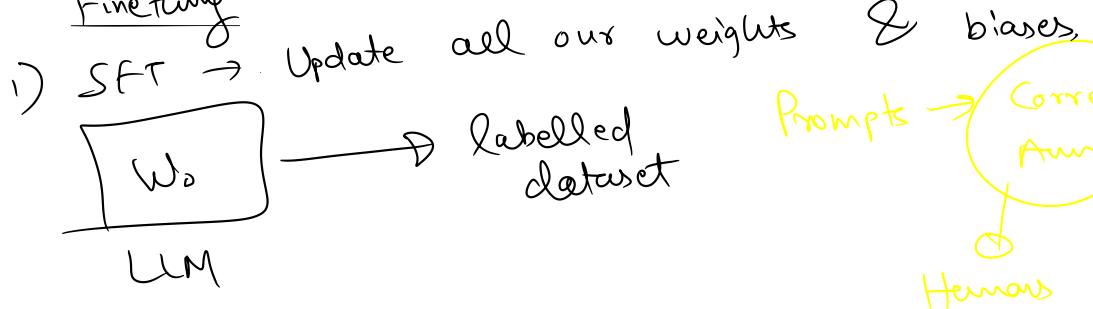
"What is today's weather in Kampur"

14th Jan

3) Hallucinations



Fine-tuning



2) PEFT : LoRA & QLoRA

$$W = W_0 + \text{AB}$$

Low-ranked matrices for fine-tuning.

Freezed

✓ 1) Private data : Amazon
Documents, Files, legal doc of Amaz

Var ~

✓ 2) knowledge Cut-off date :
Regularly update our LLM.

✓ 3) Hallucination :
"No, they never rule"

India's history
100 BCE - 2026

Dynasty, Ruler, E

Fine-tuning

Challenge

- SFT : we are updating all of our weights
- 1) Computational Cost: LORA & QLORA (A & B)
 - 2) Pretraining: English Literature. LLM French Novel
Fine tune :
 - 3) Technical Expertise: Labeled data : Frozen LLM → +ve sentiment.

"Language Models are few shot learners"

Complex P → emergent properties

IN CONTEXT LEARNING

It is the core capability of LLMs, where the model learns (GPT-3/4, Claude, Llama) to solve a task purely by seeing examples in the prompt — without updating the weight.

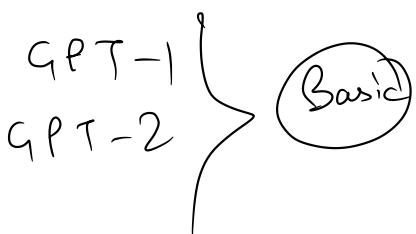
Ex

Query: Do sentiment Analysis on the following CSV file.
Please go through the below examples & learn how to judge sentiments

Examples

- 1) I love my singing → Positive Sentiment
- 2) I hate my college → Negative Sentiment
- 3) My life sucks → Negative

CSV contains! "I enjoy fishing." → ??



GPT-3 ↳

few Shot Lm
One Shot Prompt Hyp

Emergent Property: it is a behaviour or ability that suddenly appears in a system when it reaches a certain scale of complexity, even though it was not explicitly programmed on that task.

Context

Query

Query

Prompt

LM

External Database

Response

Query + Relevant Context

14 Videos

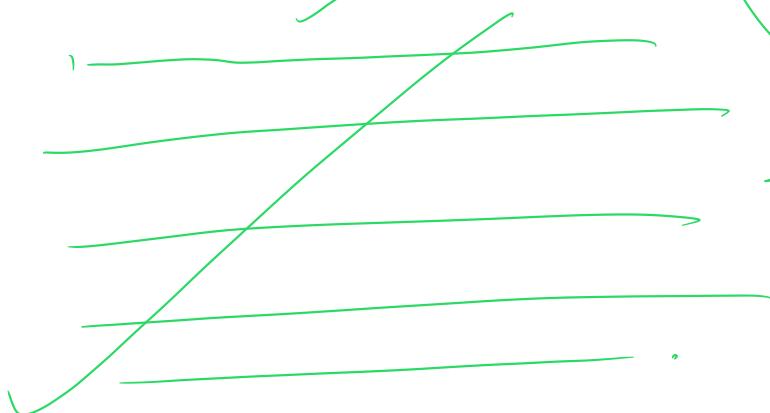
1 video

Youtube

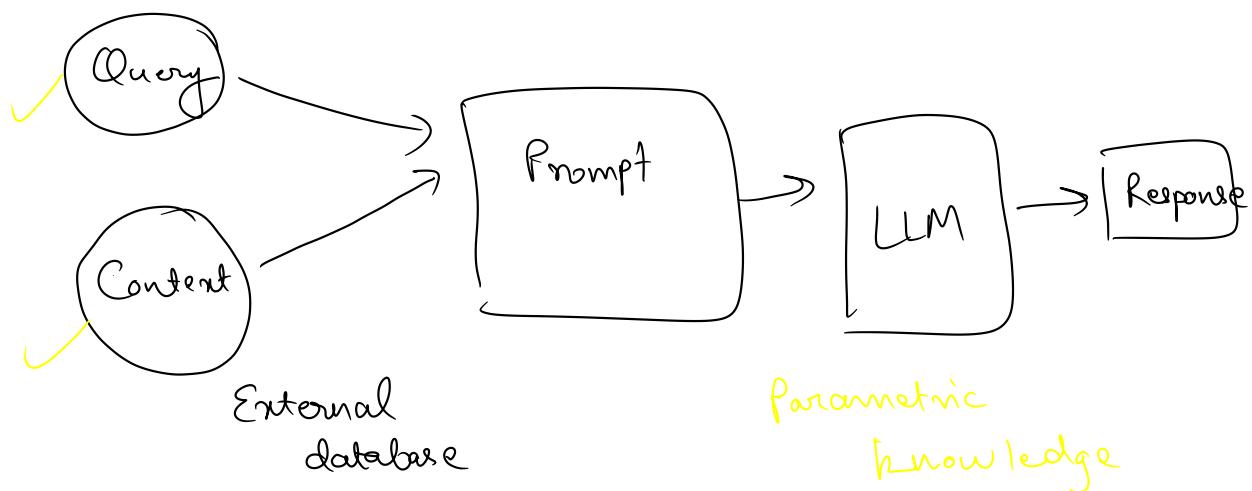
"My Tensors"

→ Response

What is Bahdanau Attention?



RAG (Retrieval Augmented Generation)

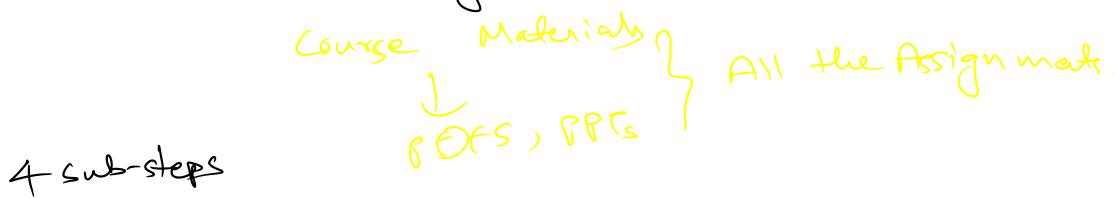


RAG: Retrieval + Generation

Indexing

It is the process of preparing our knowledge base so that it can be efficiently searched up, at the time of query

Architecting Intelligence



1) Document Ingestion:

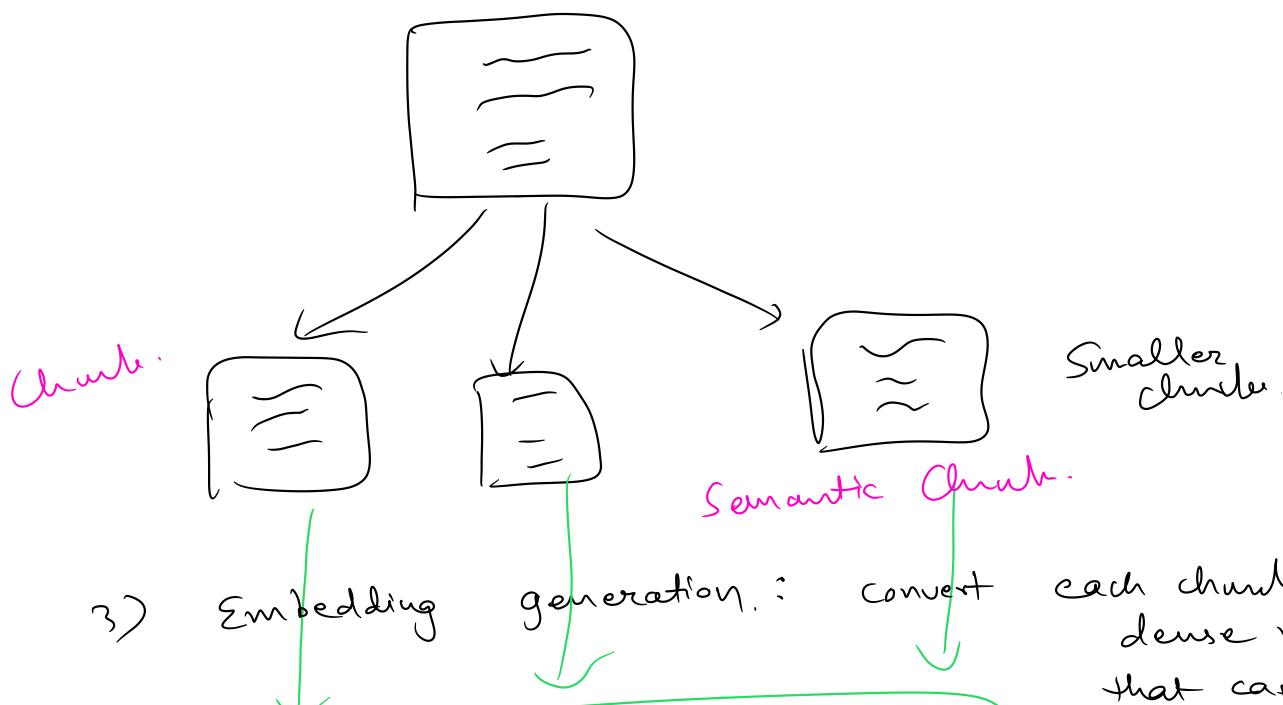
You load your source knowledge into the memory.

PDF docs
Youtube Transcripts
Github Repos

LangChain! PyPDF Loader, GitLoader etc. ←

2) Text chunking:

Break our large documents into small, semantically meaningful chunks.

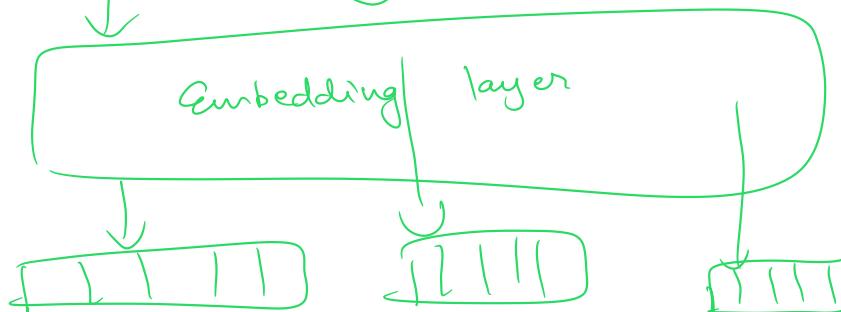


3) Embedding generation:

convert

each chunk into a dense vector that captures its semantic meaning

512 dim



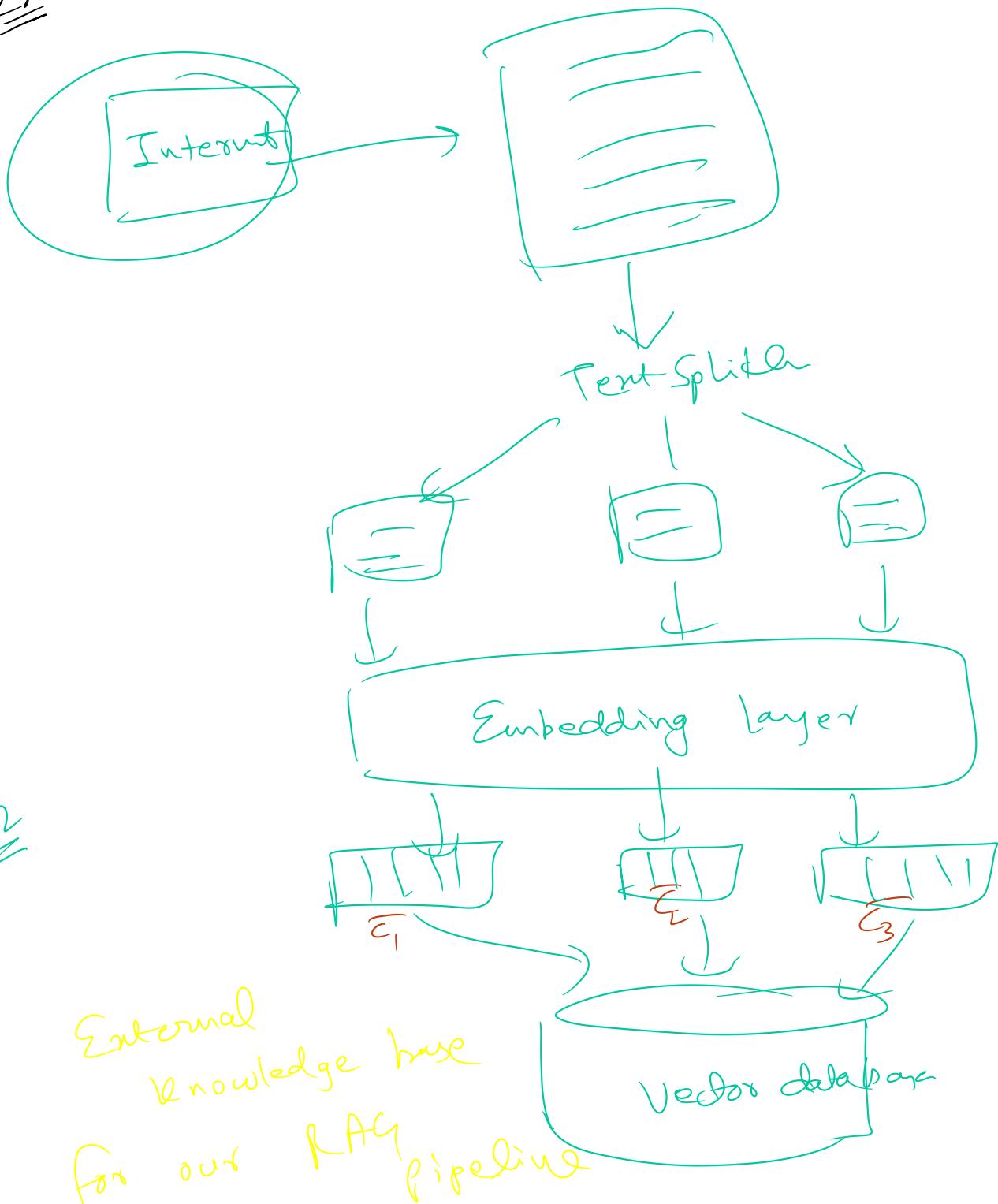
Ex: OpenAI Embedding
Gemini GenAI

4) Store : Original chunk + vector embedding

Vector database.

Chroma, FAISS, etc.

Ex

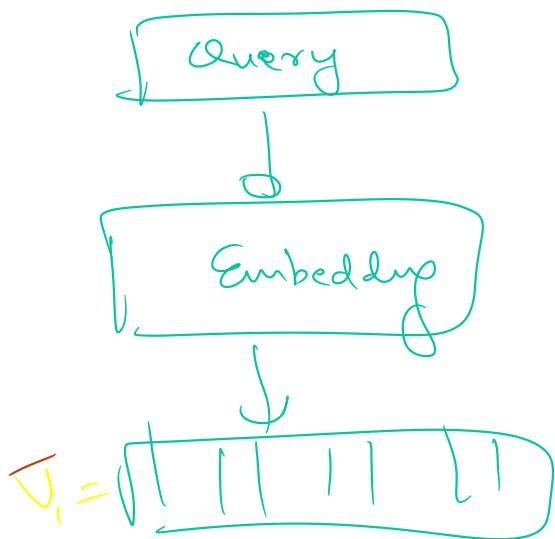


2

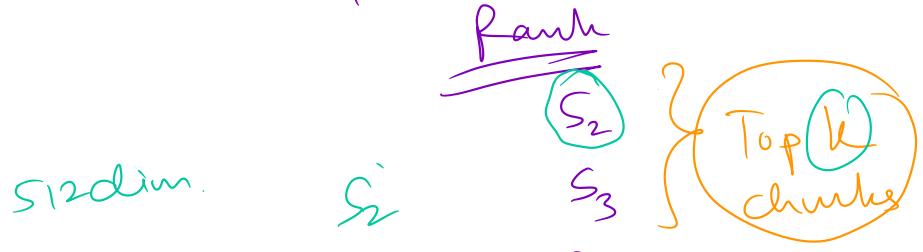
Retrieval:

finding the most relevant pieces of information from a pre-built index

Cosine Similarity, Euclidean distance, vector database
BM25



$$\text{Similarity: } \frac{\vec{V}_1 \cdot \vec{C}_1}{\|\vec{V}_1\| \|\vec{C}_1\|}, \frac{\vec{V}_2 \cdot \vec{C}_1}{\|\vec{V}_2\| \|\vec{C}_1\|}, \frac{\vec{V}_3 \cdot \vec{C}_1}{\|\vec{V}_3\| \|\vec{C}_1\|}$$



Step-1 : embedding of our query vector

Step-2 : Semantic Search

Step-3 : Ranking (Top-k)

Step-4 : Most relevant

3

Augmentation

Where the retrieval documents (chunks of relevant context) are combined with the user's query to form a new, enriched prompt for the LM.

Example

" " You are an helpful assistant

Answer the question ONLY from the provided context. If the context is insufficient, just say you don't know.

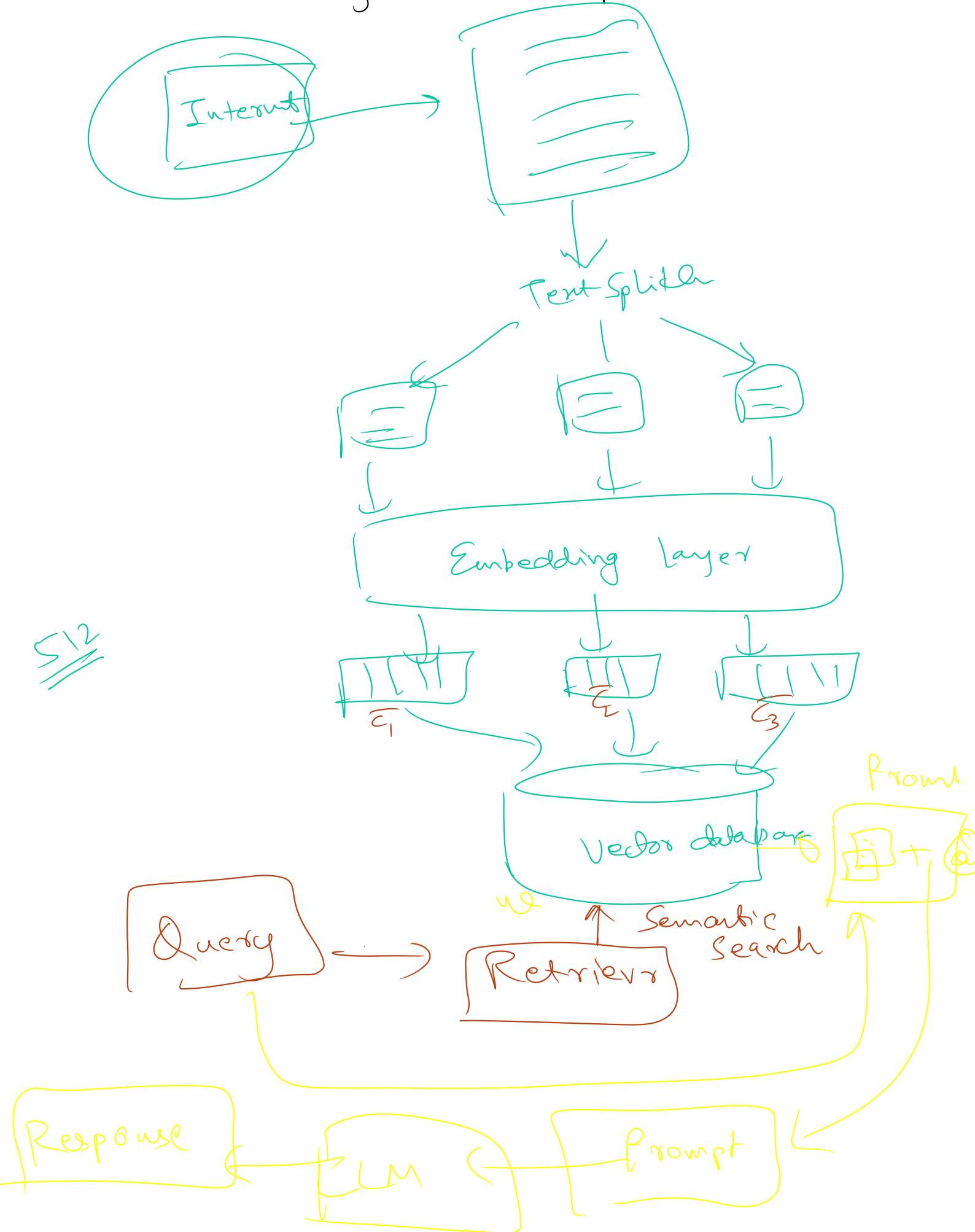
{content}

Question, {question}

" " "

④ Generation

LLM uses the user's query & retrieved & augmented context to generate a response



Vanilla LLM

RAG

1) private data

"Amazon"
"Text Ingestion"



2) knowledge Cut Off date

LLM : 1st Jan 2026

14th Jan 2026

Weather

3) Hallucinations

"Teddy Bear"

No they never ruled us"

India's history

vector

knowledge base

Legal

App RAG

Medical field: -----