# RAG

*. Retrieval Augmented Generation.
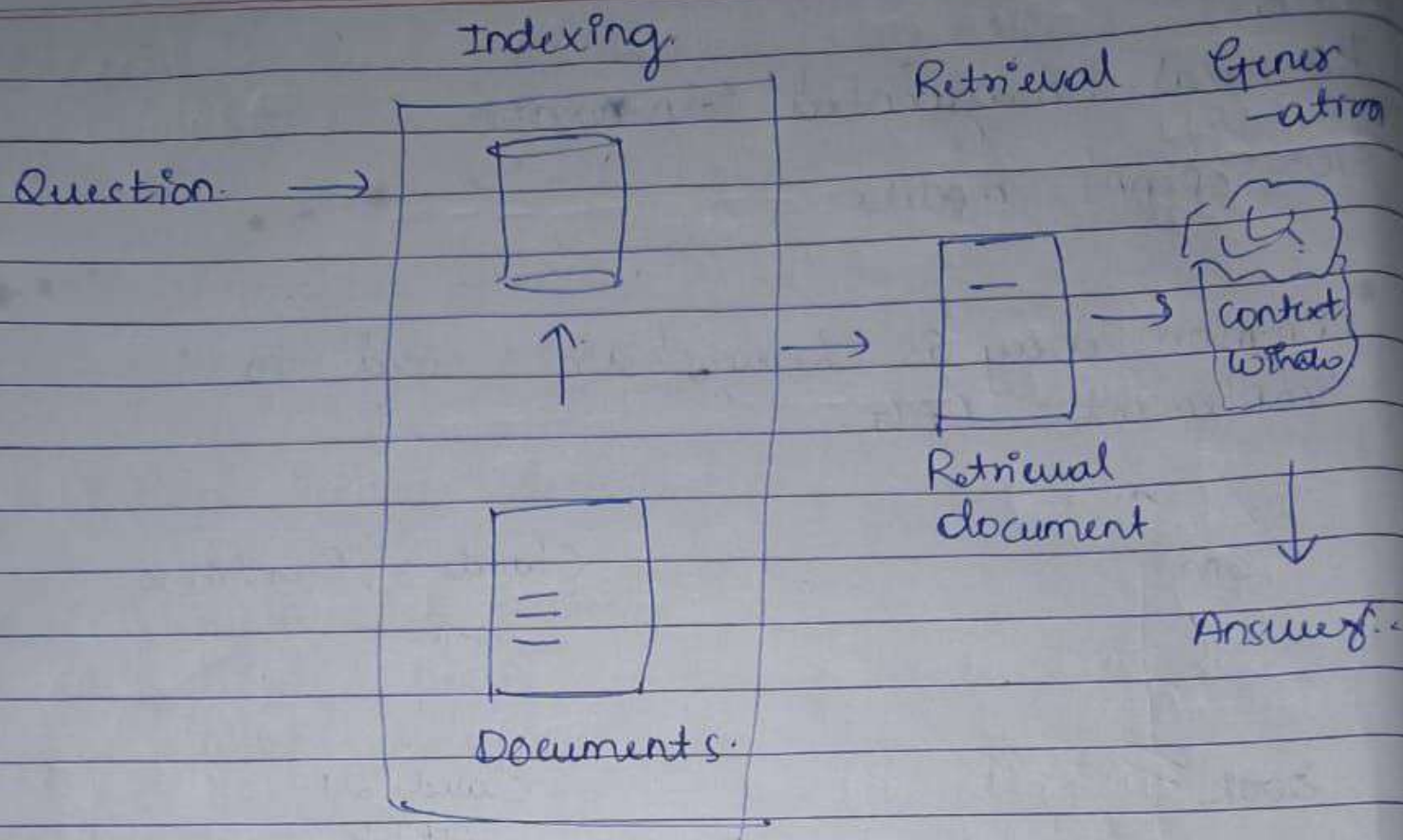
- 70% openAI models.

* common way is Langchain used to implement RAG.

*.

~1M
200k.
32k
4k

// ———→ ↓ Tokens.
                    (pretraining)

Claude 3, Gemini 1.5.
↑
↑
Claude 2.1
GPT 4

* Connecting LLMs to External data is a central need

Audio/video
I/o devices
↑↓

Tools calculator ⇄ LLM (CPU) ⇄ Search Ethernet

Data storage (disk) ⇄ Context window → ⇄ other LLM

Indexing

Retrieval    Gener
            -ation

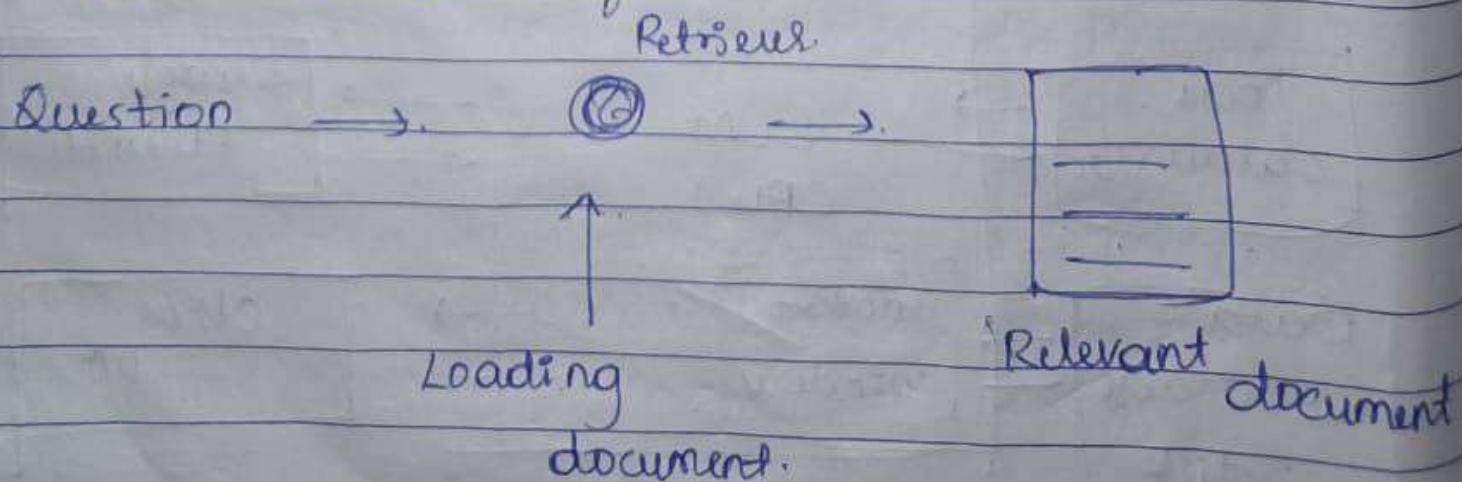Question →

context
window

Retrieval
document

Answer.

Documents.

* Basics.
  ① indexing
  ② Retrieval
  ③ Generation.

Advanced

Query Transformation
Routing
Query construction
Indexing
Retrieval
Generation

* document loading.

Retriever

Question  →      ⓒ   →
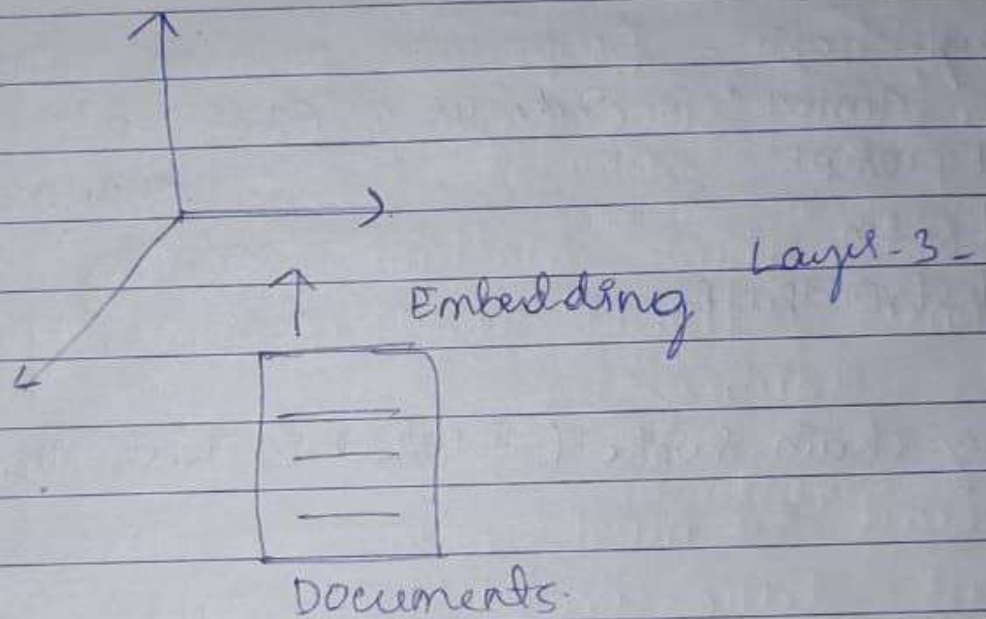
Loading
document.

Relevant
document

* Comparing Vectors.
  (Embeddings).

* Take documents
              Loading, Splitting and Embedding.

* Index makes documents Easy to retrieve
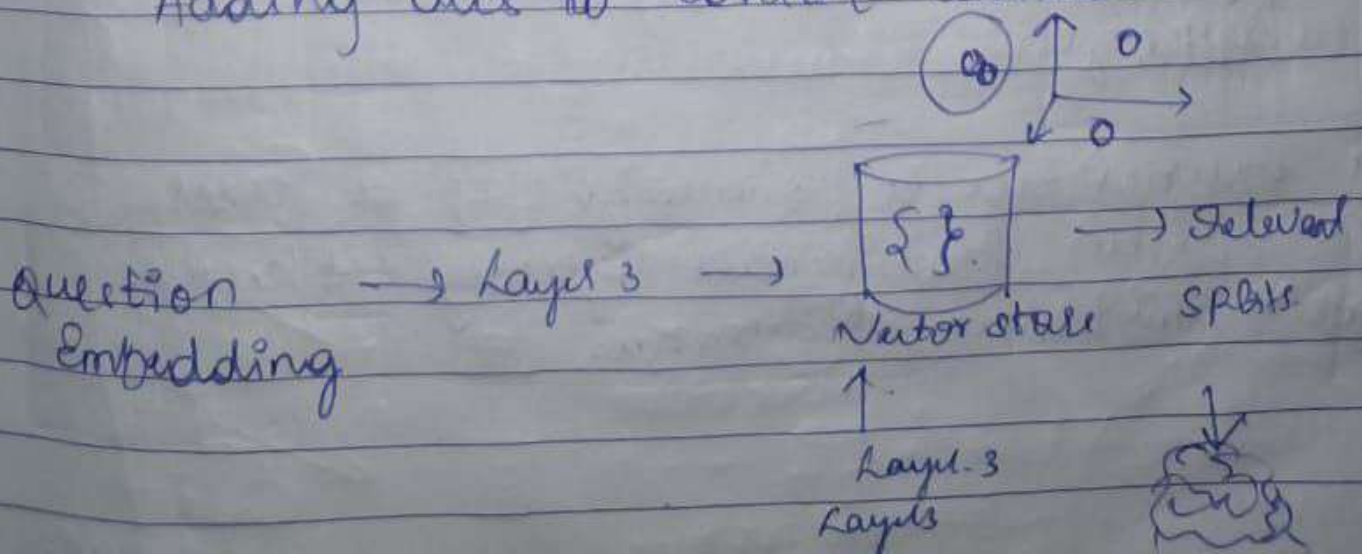
* Retrieval powered via similarity search

Embedding          Layer-3

Documents

* Langchain has Many Integration

* Generation.

       Adding docs to context window

Question → Layer 3 →
Embedding

                        Vector store      Relevant
                                          SPlits
                          ↑
                       Layer-3
                       Layers

connecting retrieval with LLMs via prompt



* rag-chain = {
  {  "context" : retriever , "question":
   | prompt                    Runnable pass through )
   | llm                                          ()}.
   | stroutputParser()
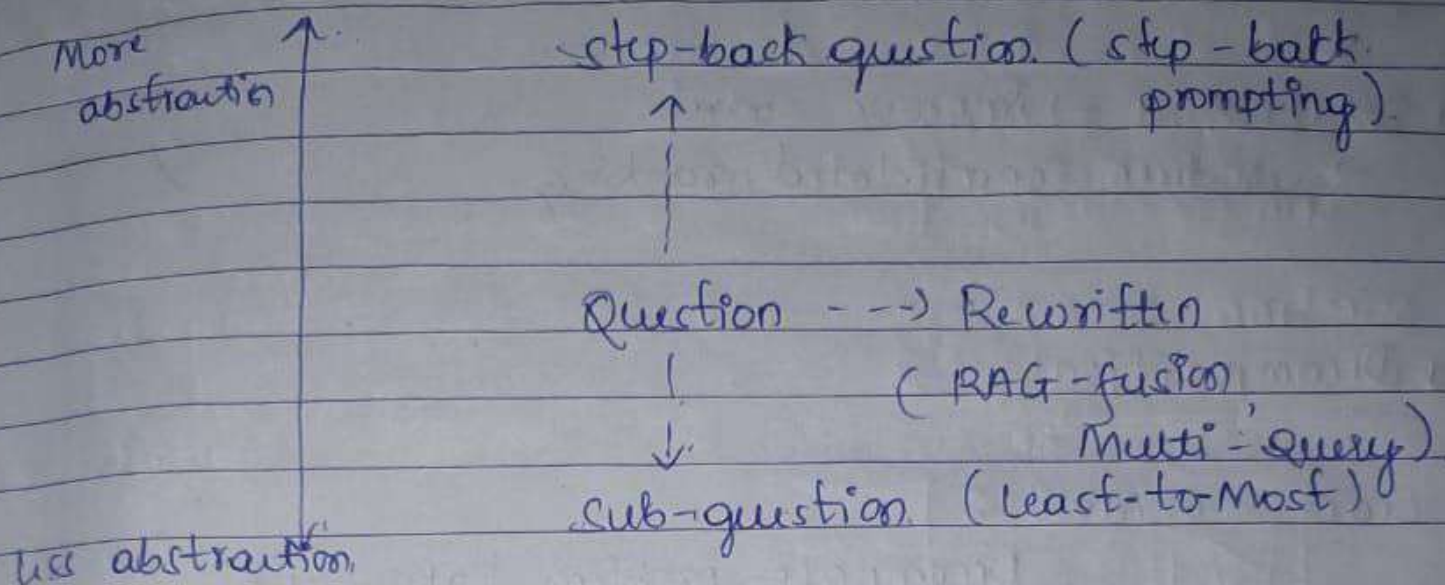  }

* rag-chain.invoke ("what is Task Decomposition")


* Query Translation
  (Multi query).

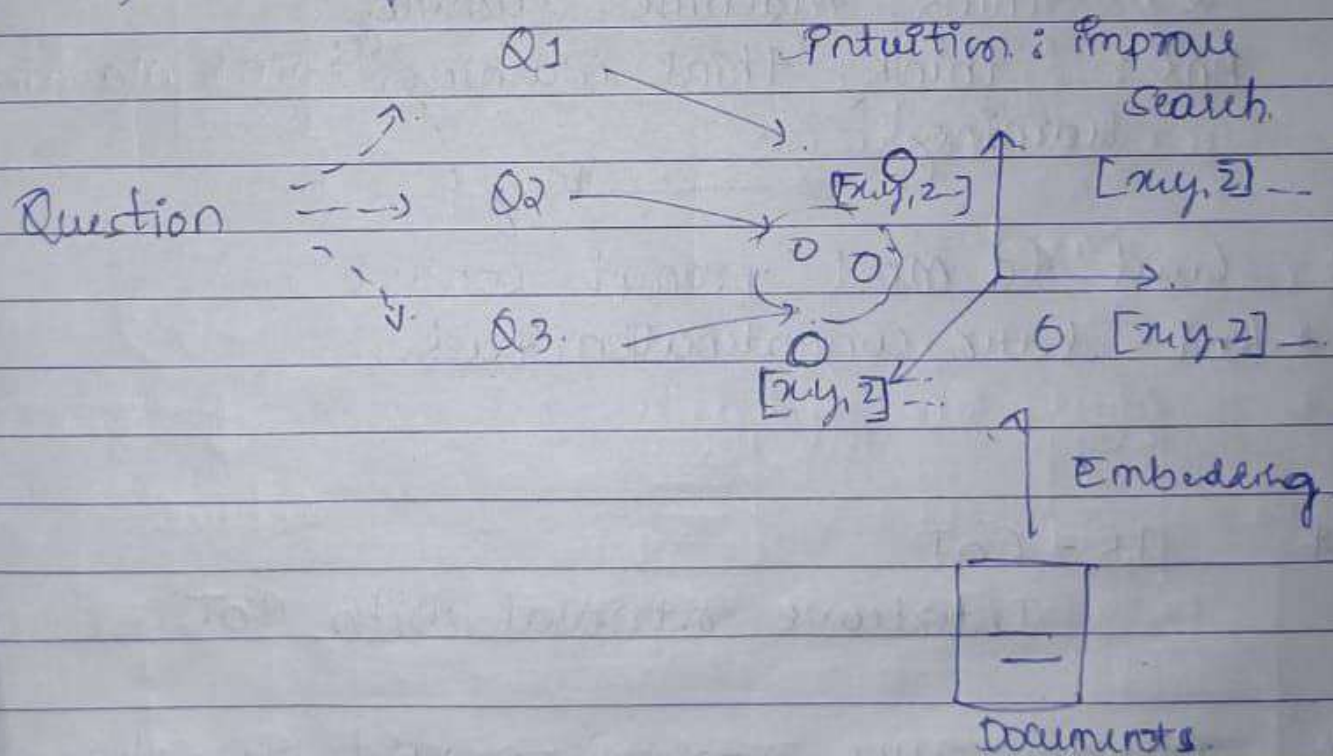Semantic search on Embeddings is hard to get right Embedding long document is a challenge.

User queries are a challenge. if a user provides an ambiguous query there get ambitious challenges
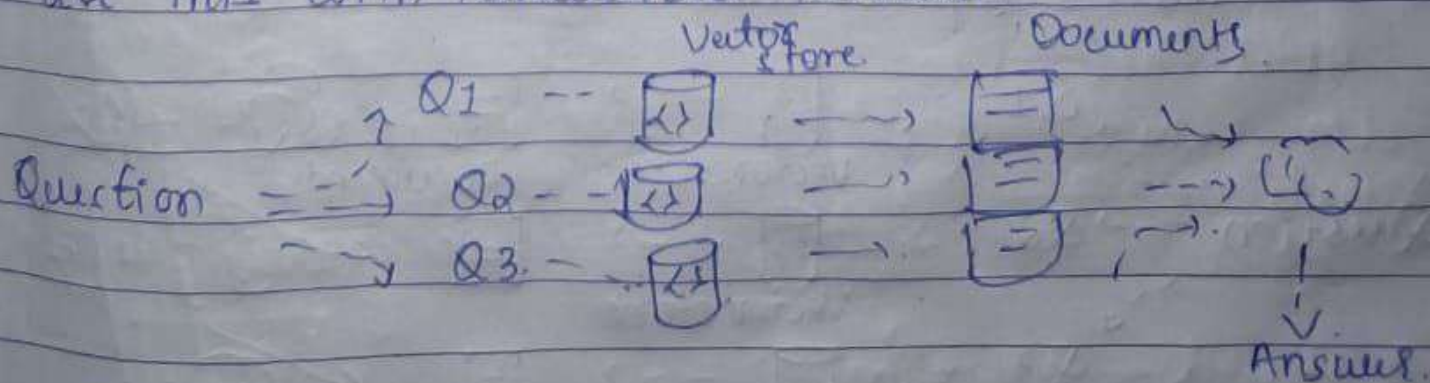
General approaches to transform questions.

More
abstraction

step-back question. (step-back
prompting)

↑

Question ---→ Rewritten
|                 (RAG-fusion,
↓                   Multi-query).
Sub-question (Least-to-Most)

less abstraction.

* Transform a question into multiple perspective.

Q1                        Intuition : improve
                                         search.
↗
Question ---→ Q2 ———→ [xy,z]↑ [xy,z] —
↘                    (° °)       →.
↓ Q3. ———→ O        6 [xy,z] —
               [xy,z]—.         ↑
                                   Embedding.

Documents.

* Use this with Parallelized Retrieval.
                    Vectorstore.    Documents.

      ↗ Q1  - -  [<>] --→ ▤
Question --→ Q2 - -[<>]  --→ ▤  --→ 🙂
      ↘ Q3 — [<>]  ——→ ▤     →.
                                    ↓
                                 Answer.

* <u>Query Translation</u>

RAG - Fusion.

* Intuition → improve search.
→ produce consolidated ranking

* <u>Decomposition</u>

Least to Most:
Decompose problem into sub problems, solve sequentially.

Q: "think, Machine, learning"
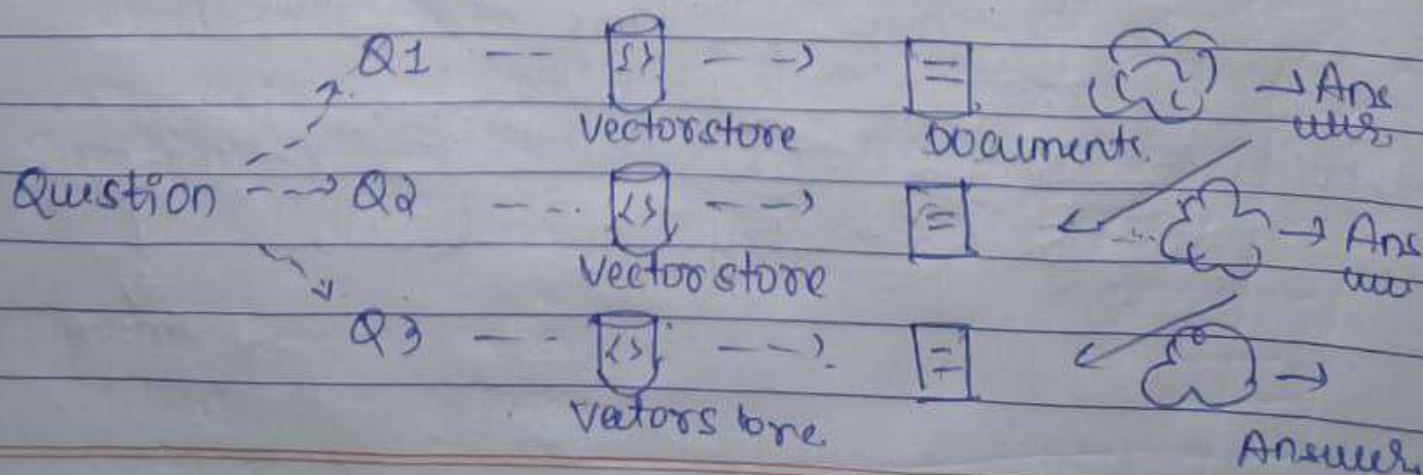Ans: ' think, "think Machine", "think machine learning'
                    Reg.

* Least to most prompt context
* last - letter concatination task.

* IR - CoT
Interleave retrieval with CoT.

* Dynamically retrieve to aid in solving the subproblems.

* Template
  (Question + Context)

* Step-back
  prompting     Independent retrieval

* Main idea is to take the question and modify
  in the way that Enhances the retriever.

* different approach suggested by Google
  where try to ask more abstract questions

* Few shot prompt to produce more abstract
  step-back question.

  * normal-context

  * stepback context.

* HyDE

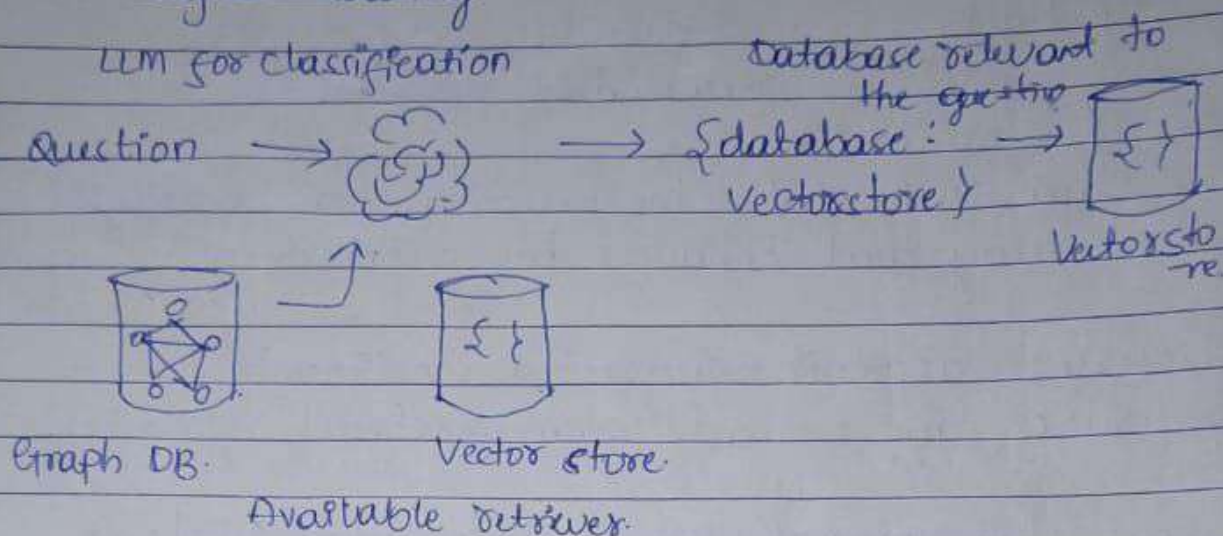$$sim(q,d) = (Enc_{(q)}(q), Ened(d)) = \langle V_q, V_d \rangle$$

* Translating raw questions into hypothetical
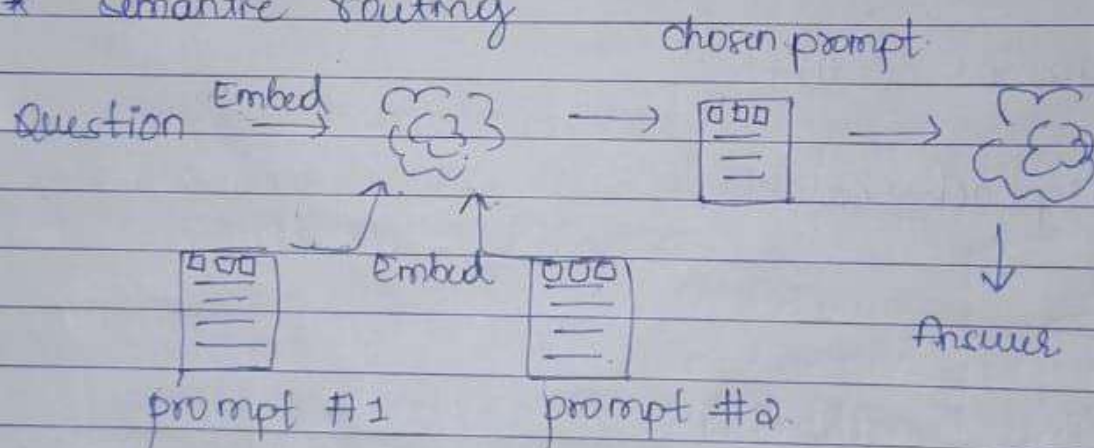  documents and that are better suited for
  retrieval

Taking question & route to the appropriate data sour

* Routing → Critical process that
(Logical + Semantic) directs user queries
to the most appropriate
* Logical Routing data source.

LLM for classification          Database relevant to
                                   the question
Question → 🌸 → {database: → [ {} ]
                  Vectorstore }        Vectorsto
                                        re.

[Graph DB illustration]          [Vector store illustration]
{}

Graph DB.              Vector store.
        Available retriever.

* Semantic routing
                              Chosen prompt
Question  Embed  🌸 → [prompt] → 🌸
                                   ↓
    [prompt]  Embed  [prompt]    Answer
    prompt #1        prompt #2.

*. LLM with structured output (object). Apply
   Bind function to     LLM returns a    parser
   LLM, LLM calls function.  JSON string
   Question → 🌸 → {'datasource':
                          'vectorstore'} ◇
                            ↑
Structured → Function
output      Schema
datasource; [...]

* Query construction          content search:                    Query

                                          chat langchain

                                    Earliest publish date
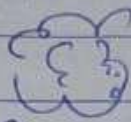
   * Query structuring          ; 2023-01-01  Apply

                          Bind function              parser

videos on                  to LLM

chat langchain  →     LLM callsfunct

published after         -ion.

2024                                    → 1. { "content search

                                              :

                            ↑            "chat langchain"

                    Earliest publish date:

         Function schema          "2024-01-10" }

Vector DB  metadata
   Content search : [...]
 Earliest publish date : [...]


* unstructured input to structured query.

   input  →        "          "

            ↓

   o/p  →.  Context search  -  "    "
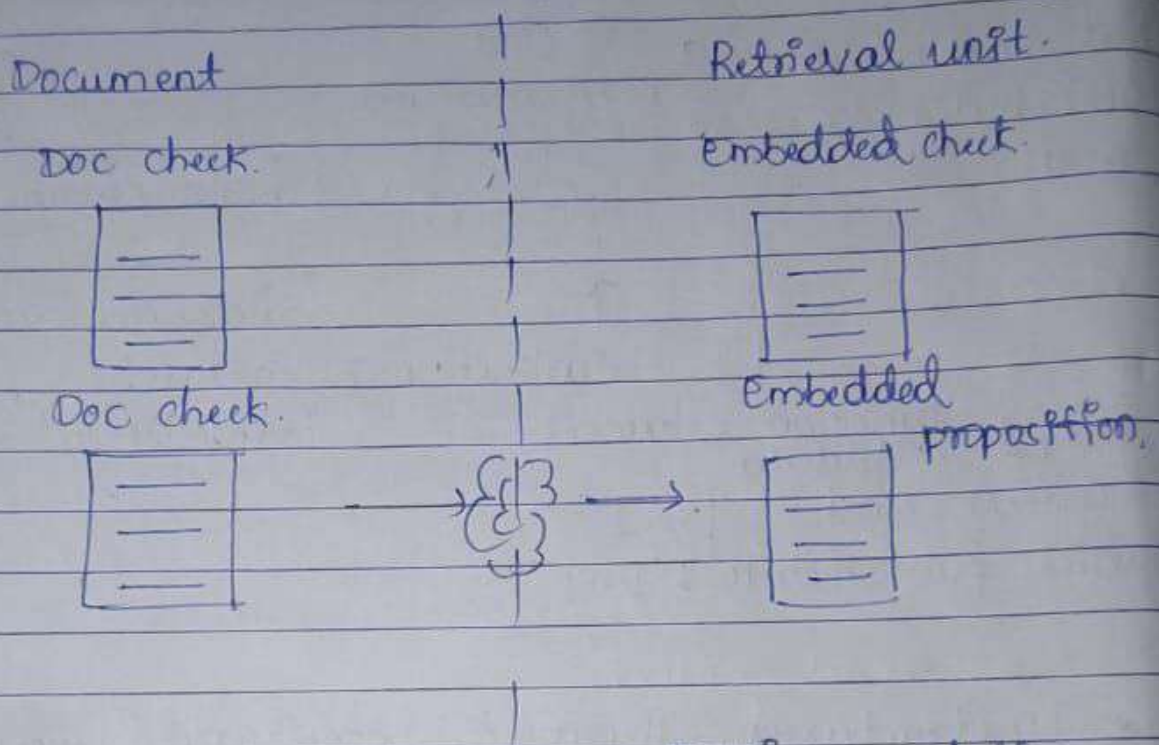
            title search   -  "    "

            Max_length-sec   300.
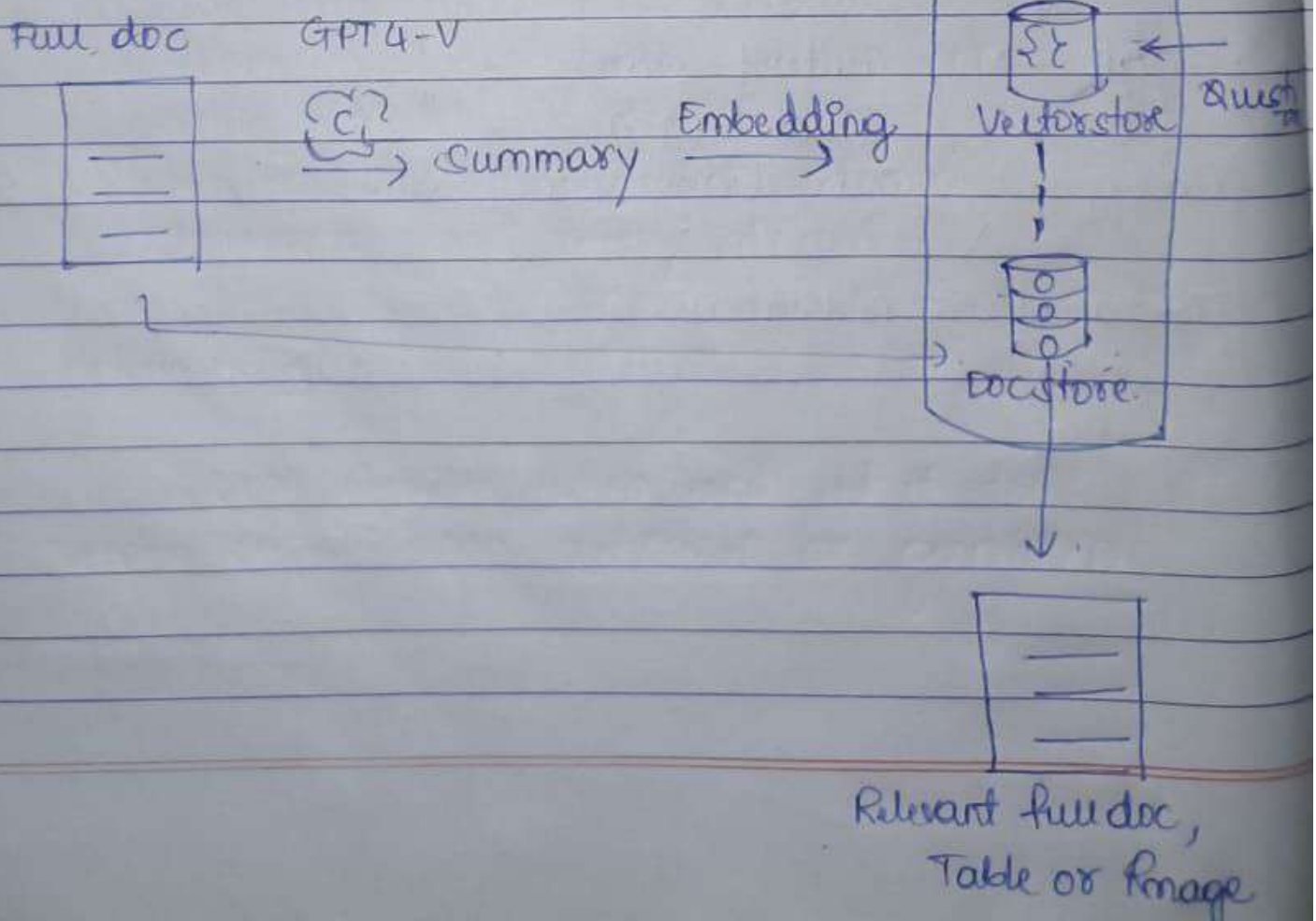

Meta data filtering

* Indexing (Multi-represcentation)

* for vector store.

Propocition Indexing.

Document | Retrieval unit.

Doc check. | Embedded chuck

Doc check.

Embedded propacition.

Multi representation Indexing

* Multi vector Retriever.

Full doc          GPT 4-V

Summary    Embedding    Vectorstore  Quest

Docstore.

Relevant full doc,
Table or image

# * Indexing (RAPTOR)

(raw documents).



cluster summary

Embed + cluster

Summary

Embed + clusters.

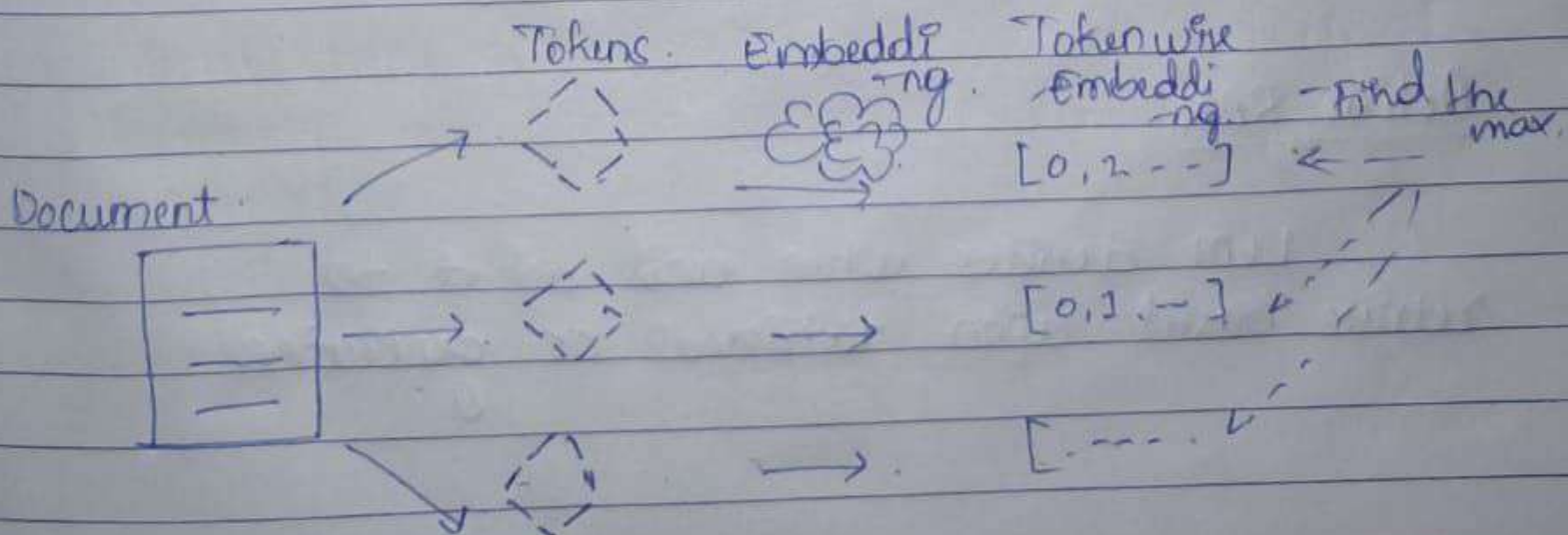summaries that integrate info across docs.

Vectorstore ↑ Summary ↕

* Build a hierarchical index. Root Summary

More abstract, high level summary.

* Indexing (COLBERT).
Specialized Embedding.

Tokens.   Embeddi-ng.   Token wise Embeddi-ng.   — Find the max.

Document

$[0, 2 --]$ ← —

$[0, 1 --]$

$[----]$

Tokenwise
Embedding                    Embedding    Tokens

[0,1 --]                         ←           ◇  ←

[0.1, 0.2 --]                    ←           ◇  ← Question

[0.1, 0.2 --]                    ←           ◇  ✓

                    Doc Score.

        Sum of max similarity of
        Each query Embedding to
        any document Embedding.


    Library  -  ragatouille.



CRAG    (Corrective RAG)


    Challenges
    ① we can ask when to Retrieve
    ② when to re-write the question for
       better Retrieval
    ③ when to discard irrelevant retrieved
       documents and re-try retrieval.


Active RAG


        LLM decides when and what to
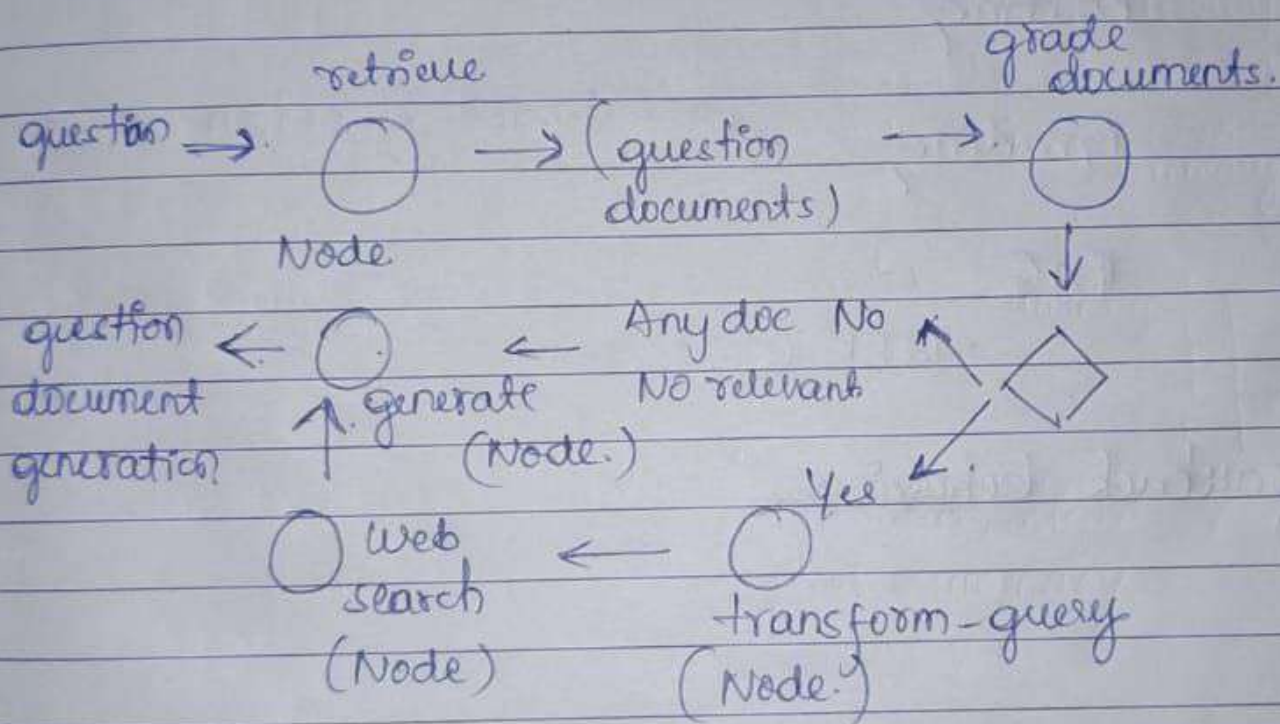    retrive based upon retrieval or generation

* <u>Levels of control</u>

Cognitive architecture.

* State machines are good way to implement active RAG
    * Knowledge Refinement stage.
    * using Langchain than the Langgraph is relevant here.

question → retrieve ⟶ (question documents) ⟶ grade documents.

Node.

question ← ○ ← Any doc No      No relevant

document generation ↑ generate (Node.)

○ Web search (Node)  ←  ○ Yes  transform-query (Node.)

* Adaptive RAG
  Query Analysis
  RAG + Self reflection
                    ┌→ related to index ~vector
                    │                    store
Question — Query Analysis → current Event
                    └→ anything else
                          chroma (vector DB).

*. Command - R
      Routing
      Grading    ⌐→ relevant question or not

      RAG
        - 128K context.

* LLM fallback behaviour.