



From Vectors to Answers: Practical RAG with Qdrant



MOHAMED ARBI NSIBI

- **ML engineer @ Carbon Insights**
- **Final year ICT engineering student@ SUP'COM**
- **Former GDSC Lead 23/24**
- **GDGoC SUP'COM & ISAMM Mentor**
- **GDG Carthage Co-organizer**



<https://huggingface.co/Goodnight7>



<https://www.linkedin.com/in/mohammed-arbi-nsibi-584a43241/>

Content



- Motivation
- RAG components
- Vector stores deep dive
- Building basic RAG pipeline
- Some Advanced RAG Techniques



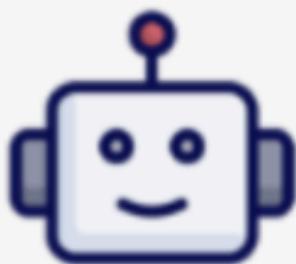
1- The need for an external Knowledge !

2- Hallucinations

Hallucinations



Is 9677 a prime number?



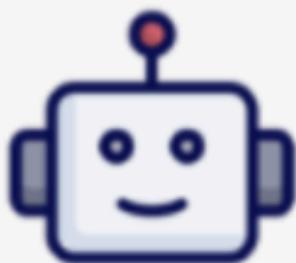
No, 9677 is not a prime number.

It can be factored into 13 and 745, as $9677 = 13 \times 745$.

} incorrect assertion
} snowballed hallucination



Is 9677 divisible by 13?



No

in a separate session,
GPT-4 recognizes its
claim as incorrect!

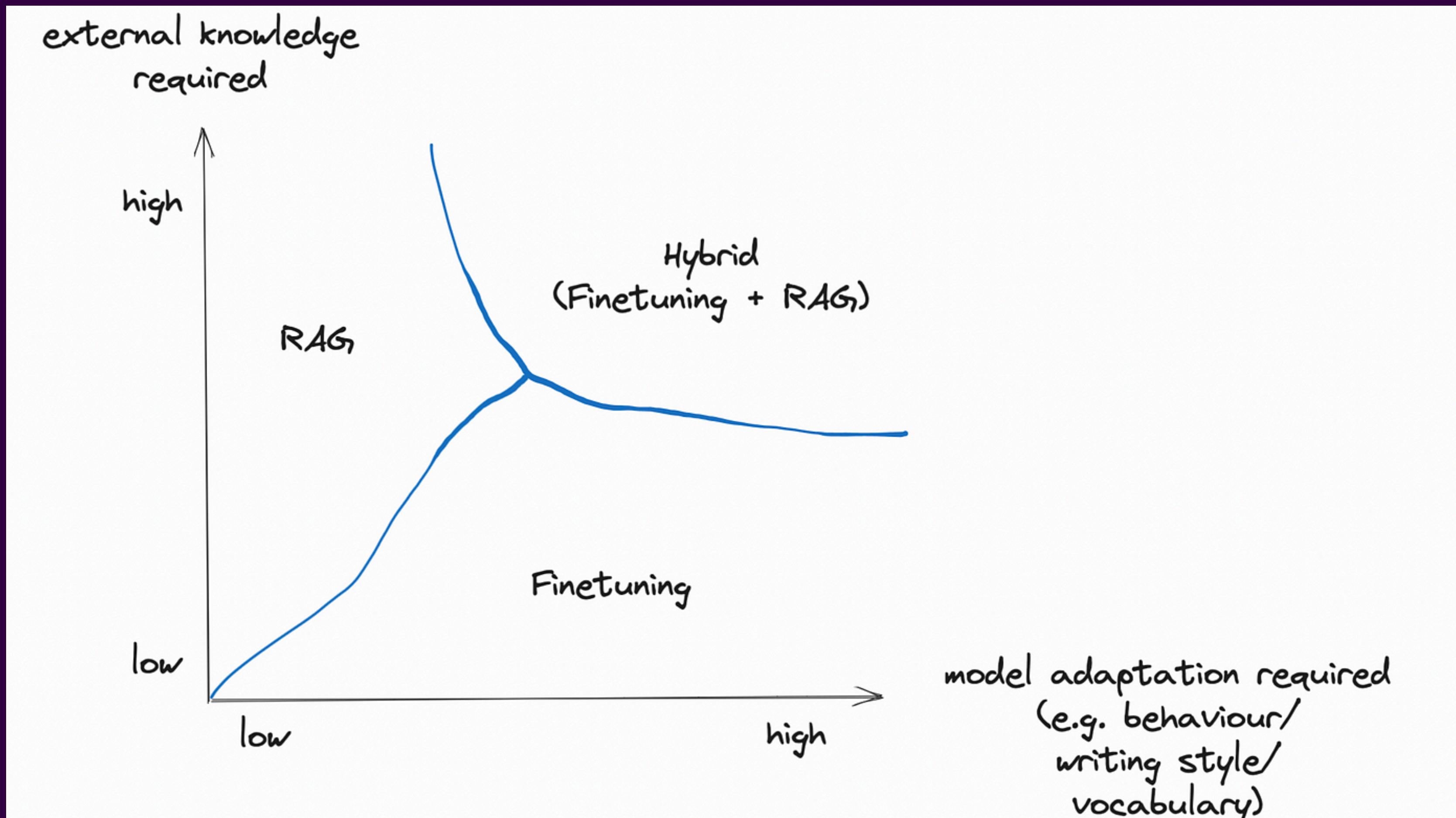


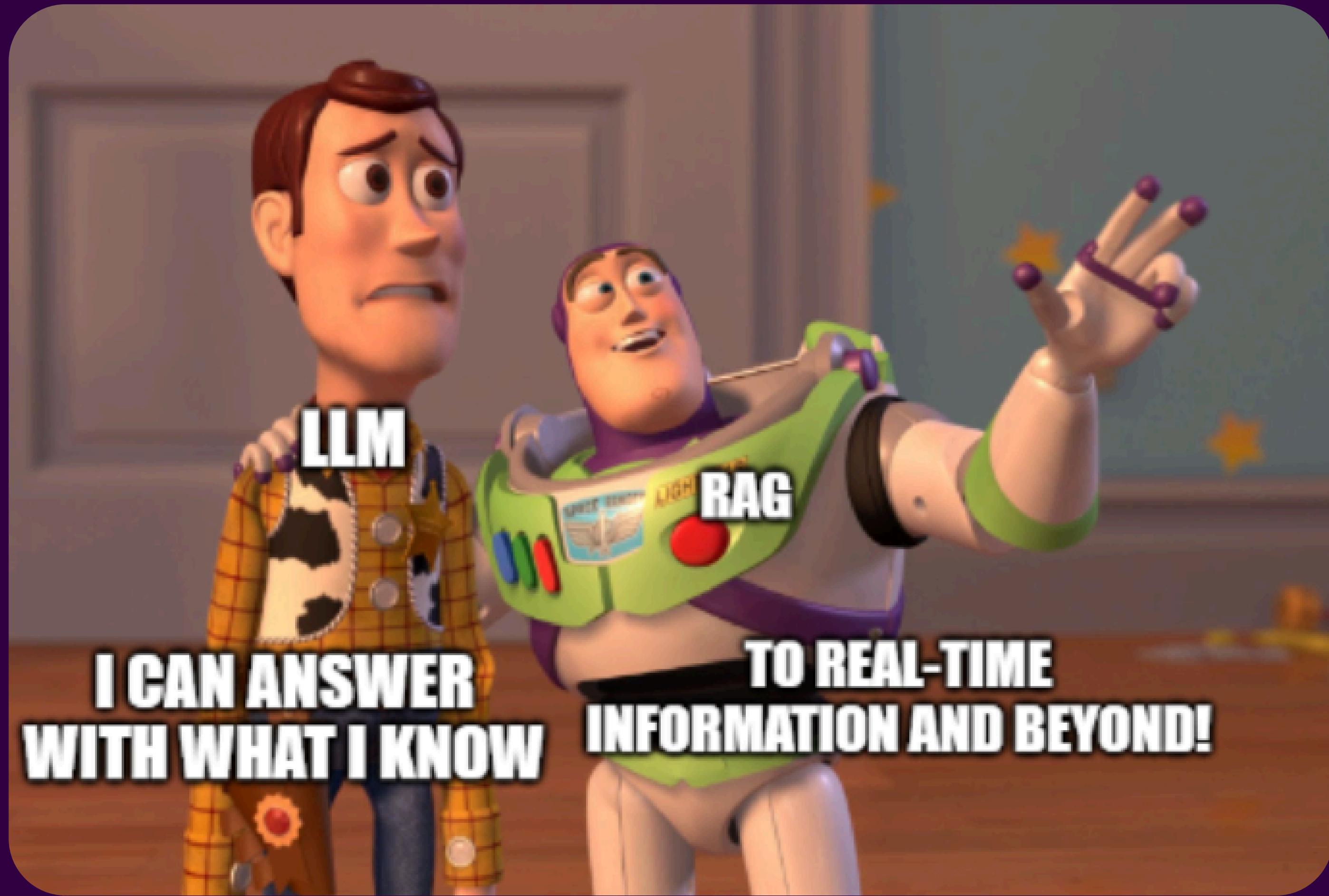
Hallucinations

- The model is not trained on enough data.
- The model is trained on noisy or dirty data.
- The model is not given enough context .
- The model is not given enough constraints (rules, guidelines, or limitations)



RAG / Fine-tuning

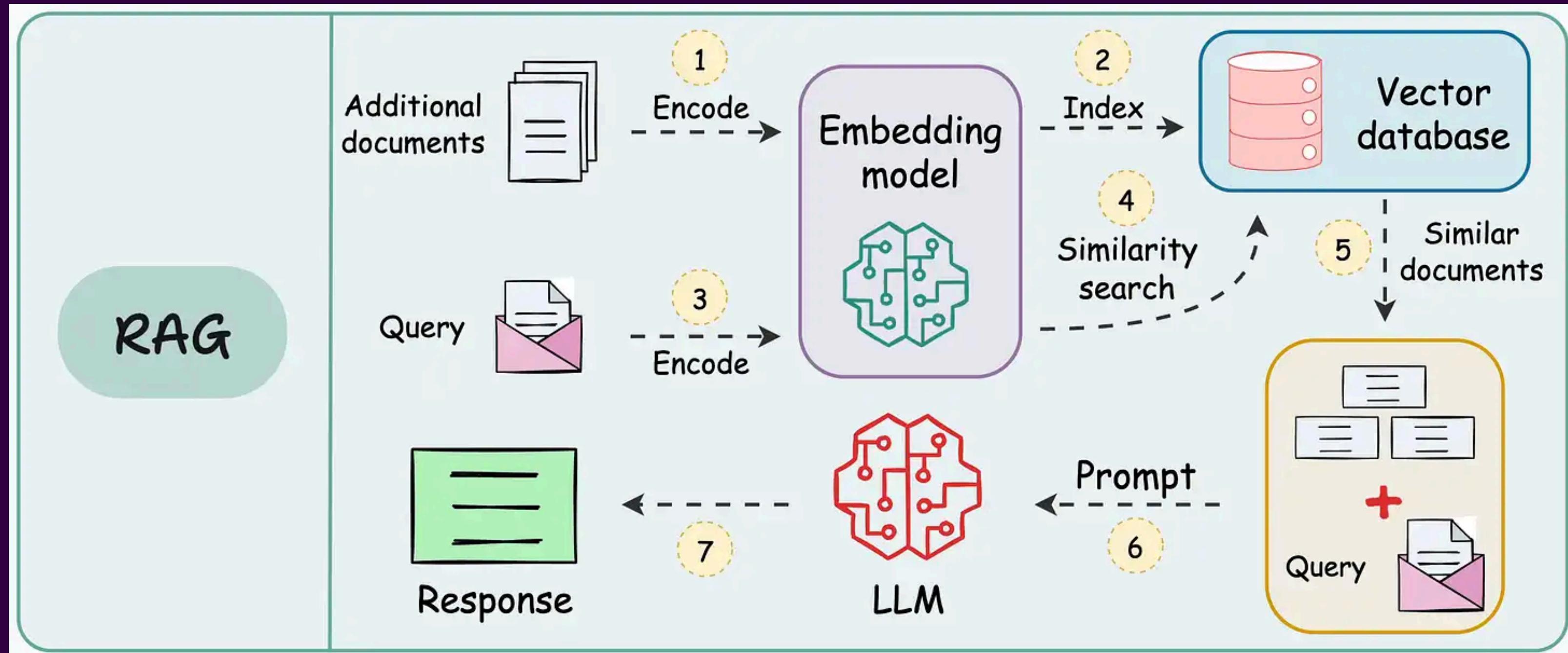




LLM
**I CAN ANSWER
WITH WHAT I KNOW**

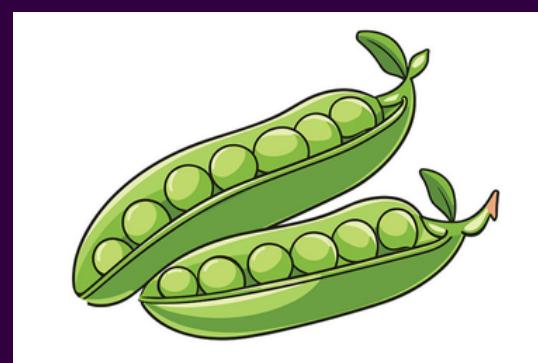
RAG
**TO REAL-TIME
INFORMATION AND BEYOND!**

RAG (Retrieval-augmented generation)





tea



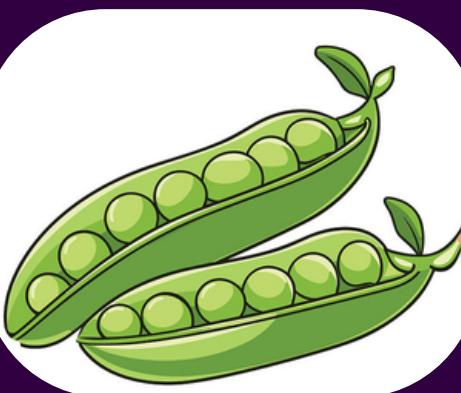
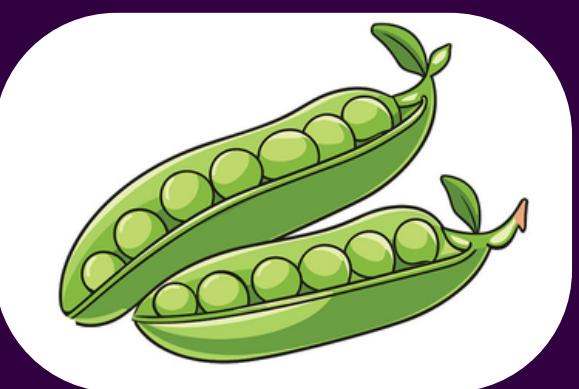
pea



coffee



≠



--	--	--	--	--



distance = 0.3

--	--	--	--	--



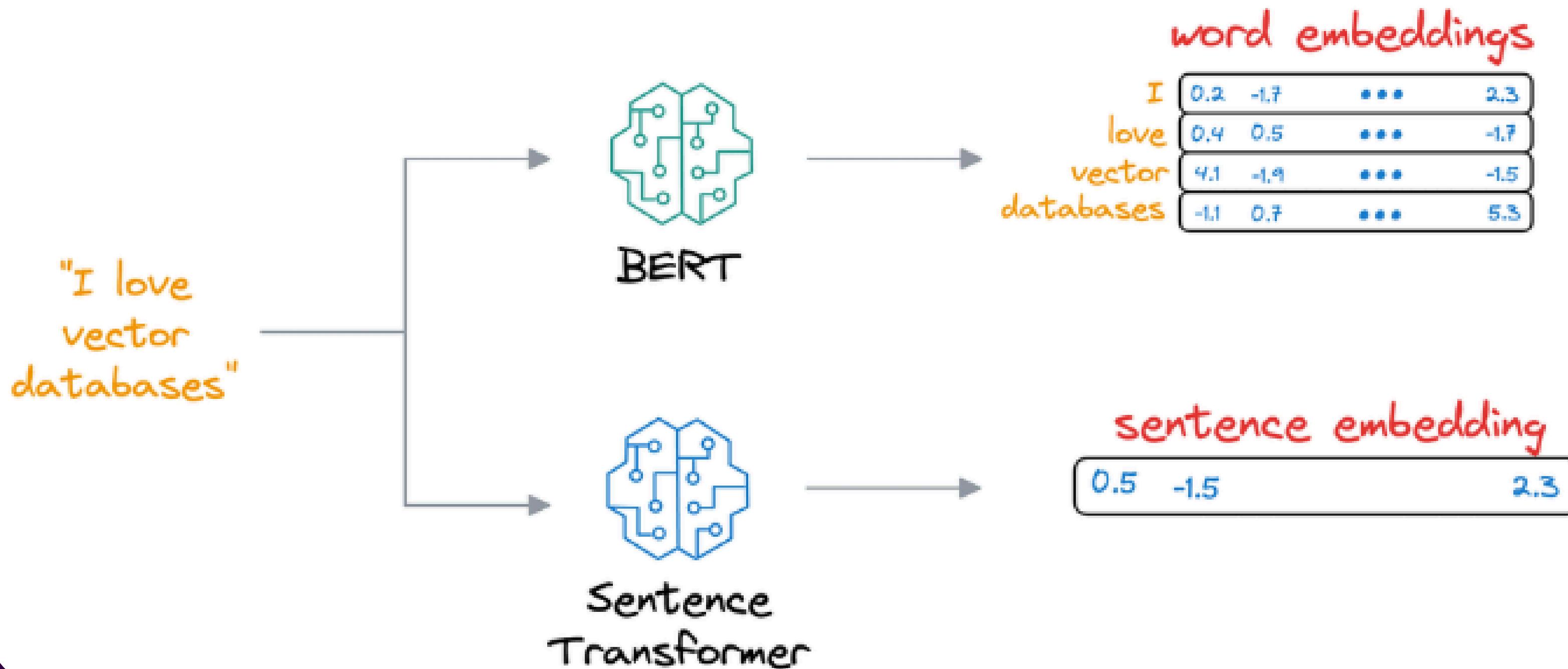
distance = 0.7

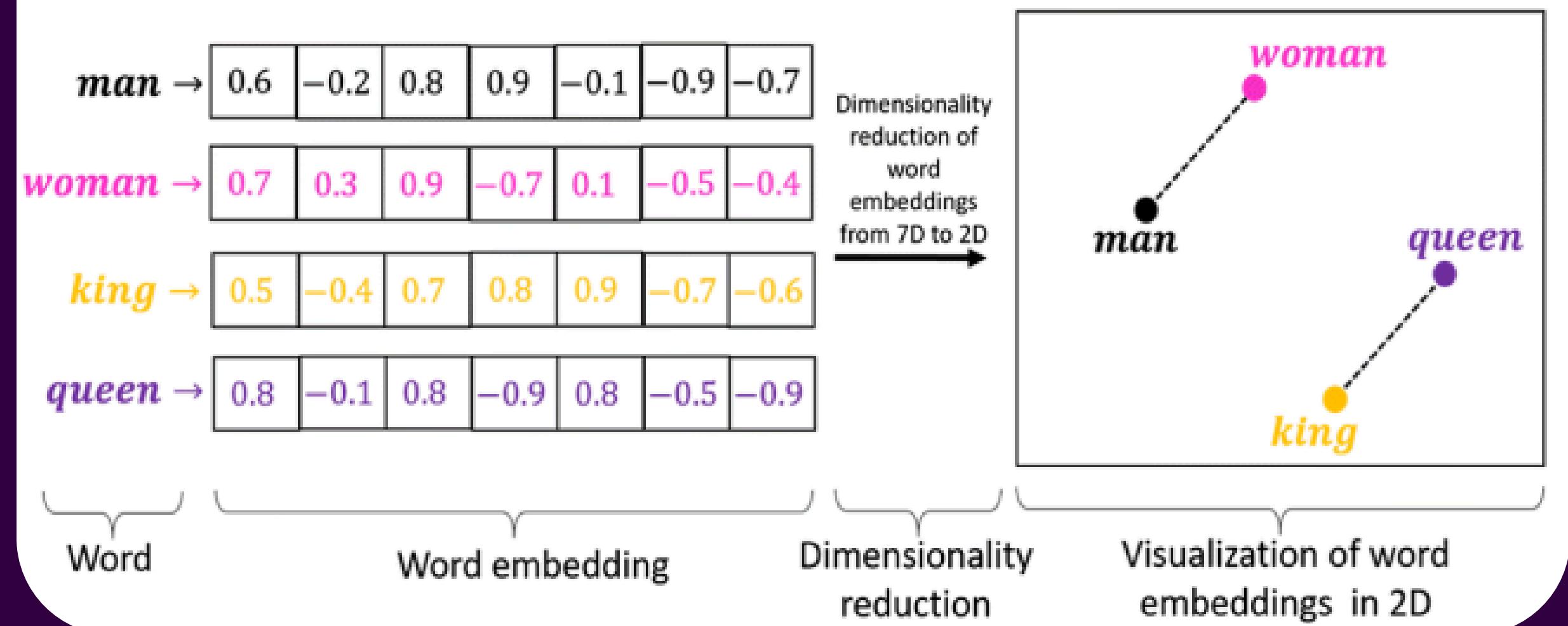
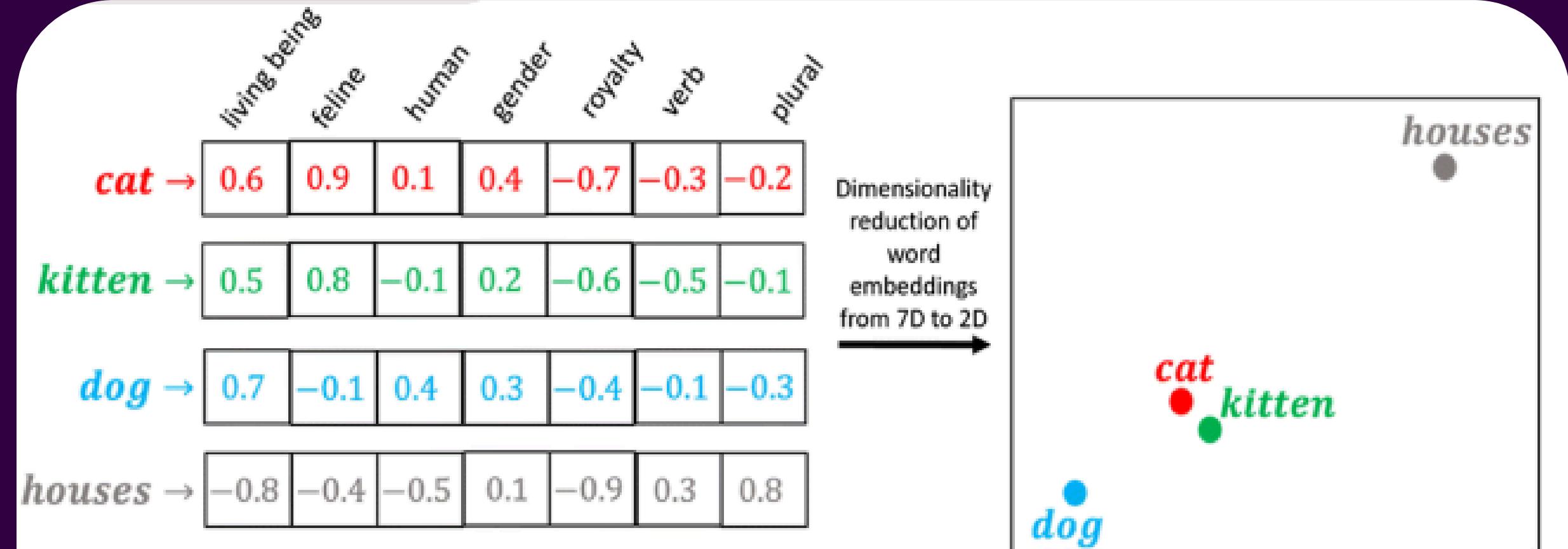
--	--	--	--	--

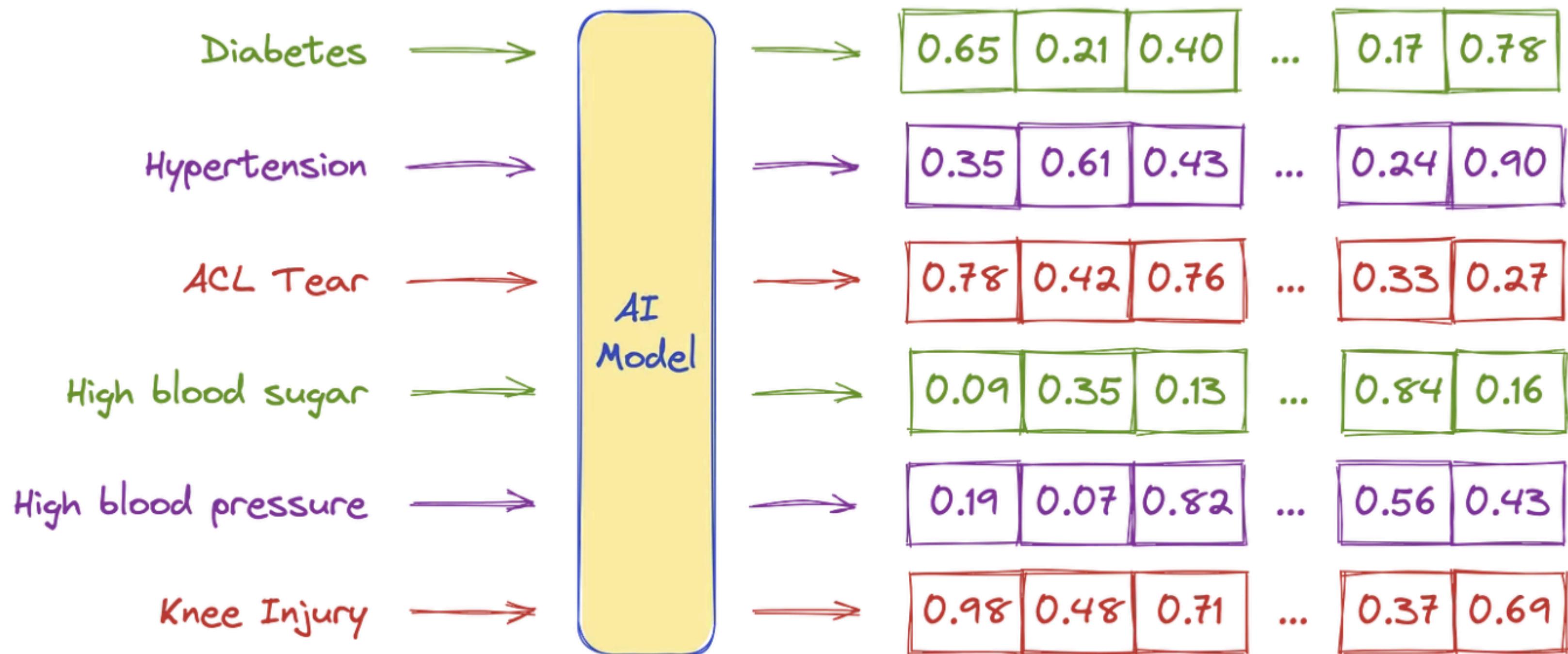
Embedding model

- These vectors live in a high-dimensional space where the proximity between vectors reflects the **relatedness** of the original items.
- Embedding model trained along LLM and learn to produce representation(vectors) based on context in word appear.

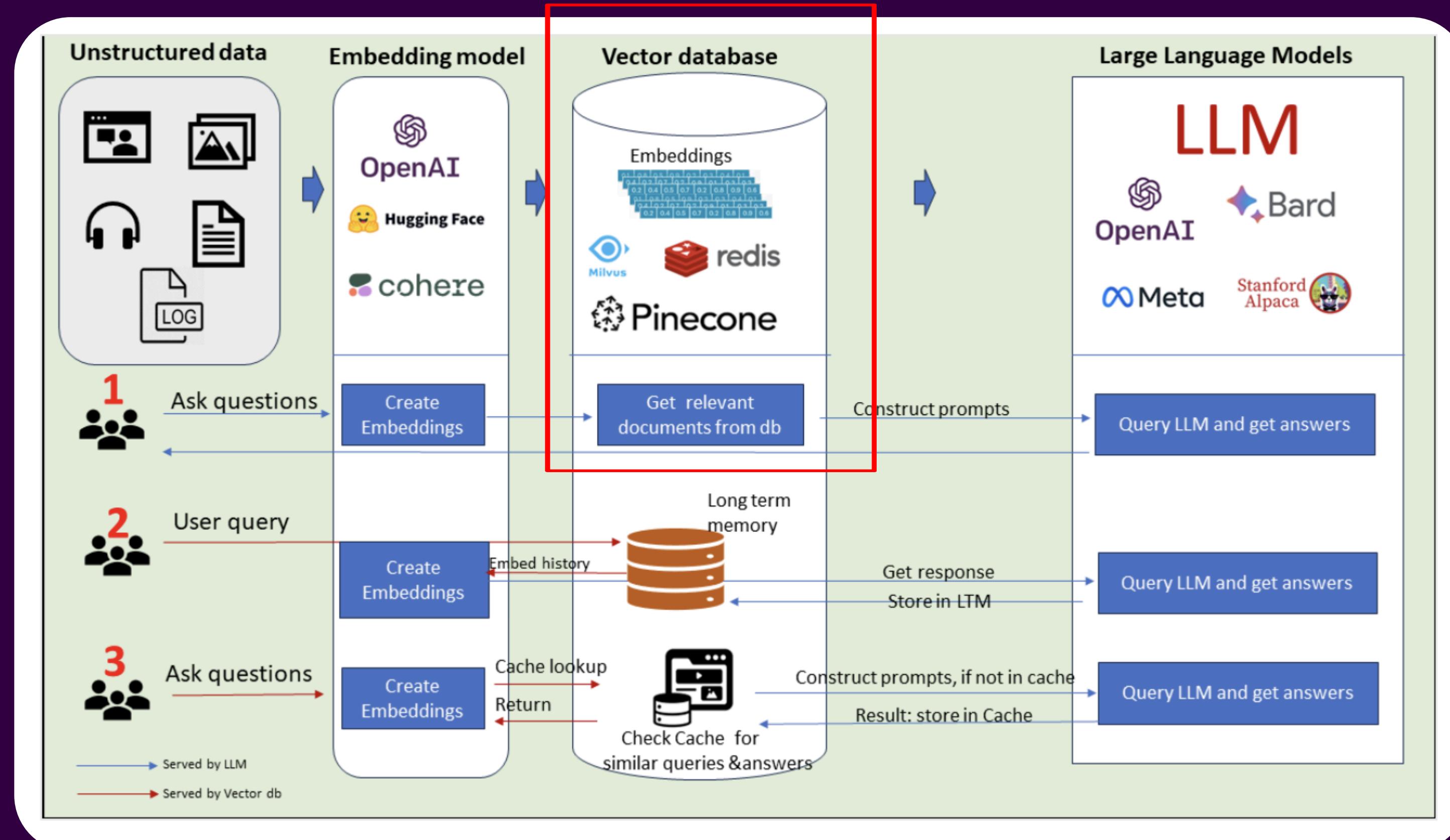
Embedding model







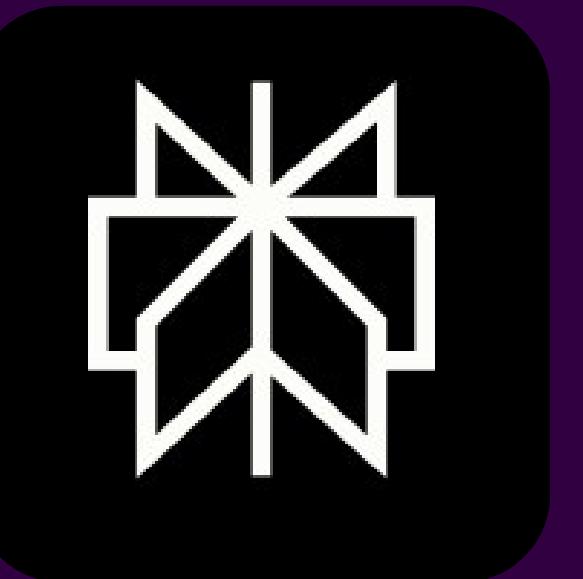
Vector Database





Fully Open-source
Self-hosting : run on ur own infra
Super fast: Latency ~0.024s (~24ms)
Hybrid Search
UI support
Free Tier ~1M(vectors) 768-dim vectors

Used by



Used by

The screenshot shows a Qdrant interface with a search bar at the top. A red circle highlights the dropdown menu under 'Database' which is set to 'Qdrant Local'. Below the search bar, there's a large 'Qdrant' logo. The main content area displays search results for repositories related to AI models. Two specific results are circled with red arrows pointing to them:

- jennifermarsman/DeepReinforcementLearning**: This repository implements the AlphaZero methodology, a deep reinforcement learning algorithm, which is a type of algorithm that can be adapted for solving complex problems like Wordle. Although it is not specifically designed for Wordle, the approach is relevant because it involves learning optimal strategies through self-play and reinforcement learning.
- jennifermarsman/PhiRecycling**: This repository uses the Phi vision model for sorting trash and recycling at scale, indicating direct use of the Phi model in its codebase.

At the bottom of the interface, there's a question from a user asking about repositories using the Phi model, and a response from a Microsoft CTO.

[NLWeb_link](#)



Andre Zayarni • 1st
Co-founder & CEO, Qdrant | Open-Source Vector Search
3w • Edited •

Kevin Scott, Microsoft's CTO, used **Qdrant** as the vector engine for the **#NLWeb** demo at the opening keynote of the annual Build Conference.

#opensource rules! 😊 🎉 🎉

You and 907 others

53 comments • 23 reposts

Love

Comment

Repost

Send



Tell them what you loved...



Most relevant ▾



Andre Zayarni Author

Co-founder & CEO, Qdrant | Open-Source Vector Search

3w ...

Recording <https://youtu.be/ceV3RsG946s?si=Xh6NkXCJH1JhKHeV&t=4879>

Like • 4 | Reply



Philippe Bourcier • 2nd

CTO #GenAI / Business Angel / Maker

3w ...

Soon a Windows build for Qdrant ?

Like | Reply • 2 replies

Why Qdrant is super fast ?

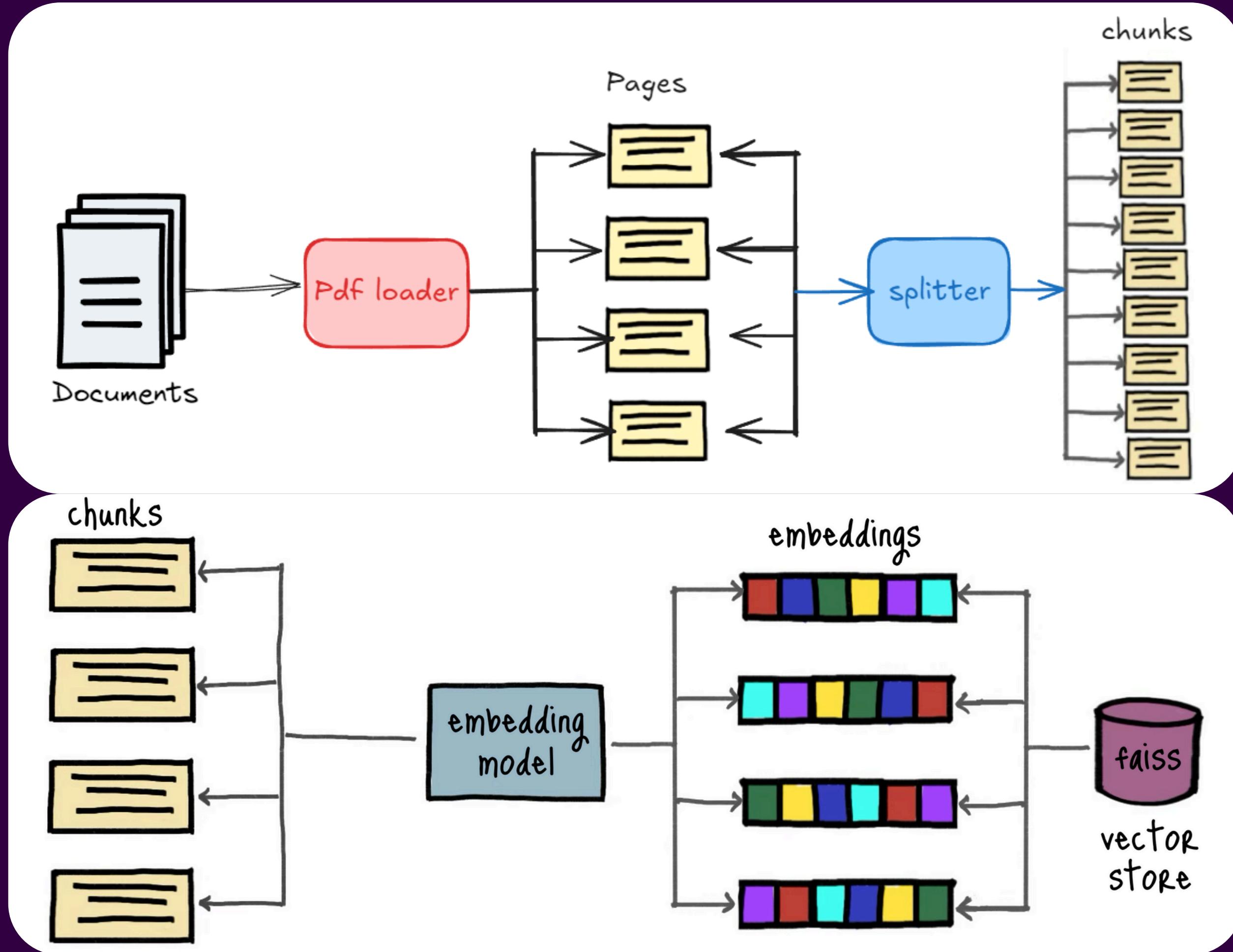


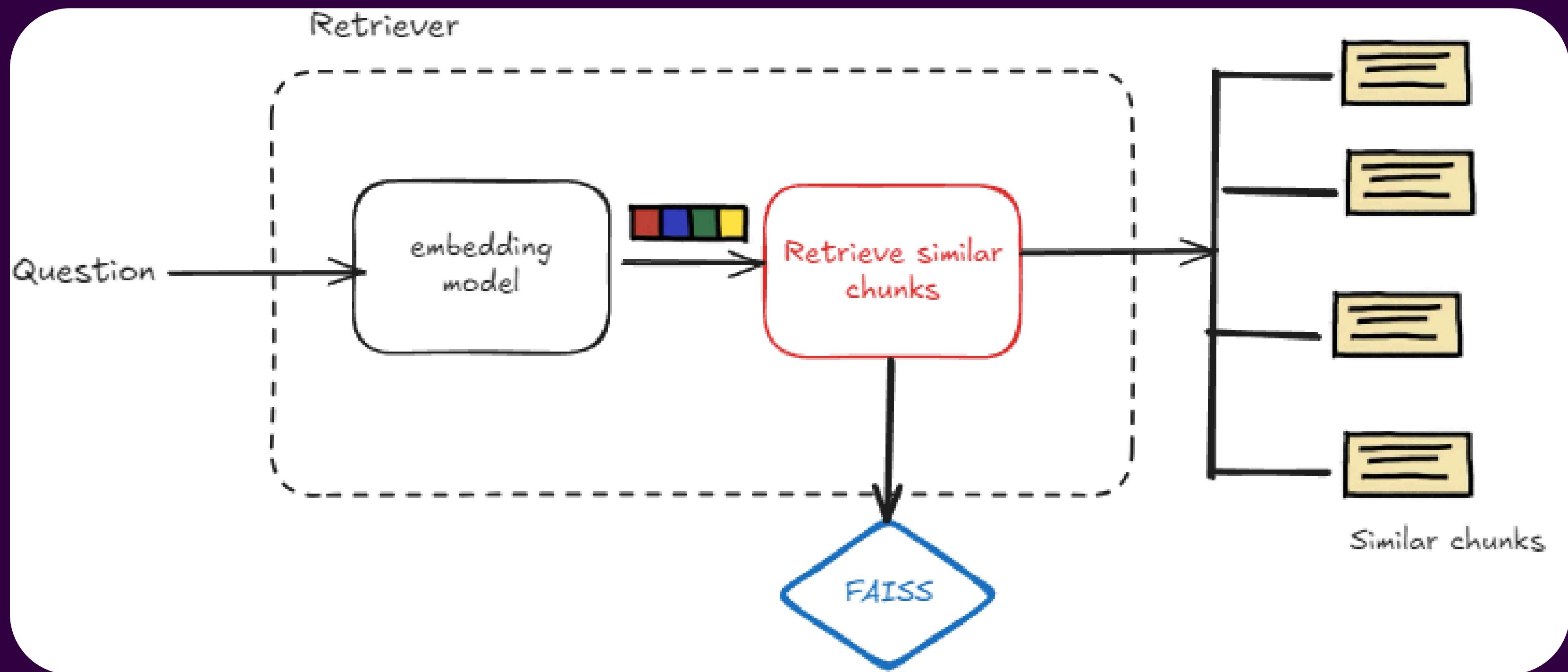
- Rust-Based Engine
- HNSW (Hierarchical Navigable Small World) Indexing : fast ANN search (on the best matches without scanning everything.)
- Vector Quantization (useful for large-scale datasets) : Saves RAM (up to 16x)
- Batch & Parallel Processing

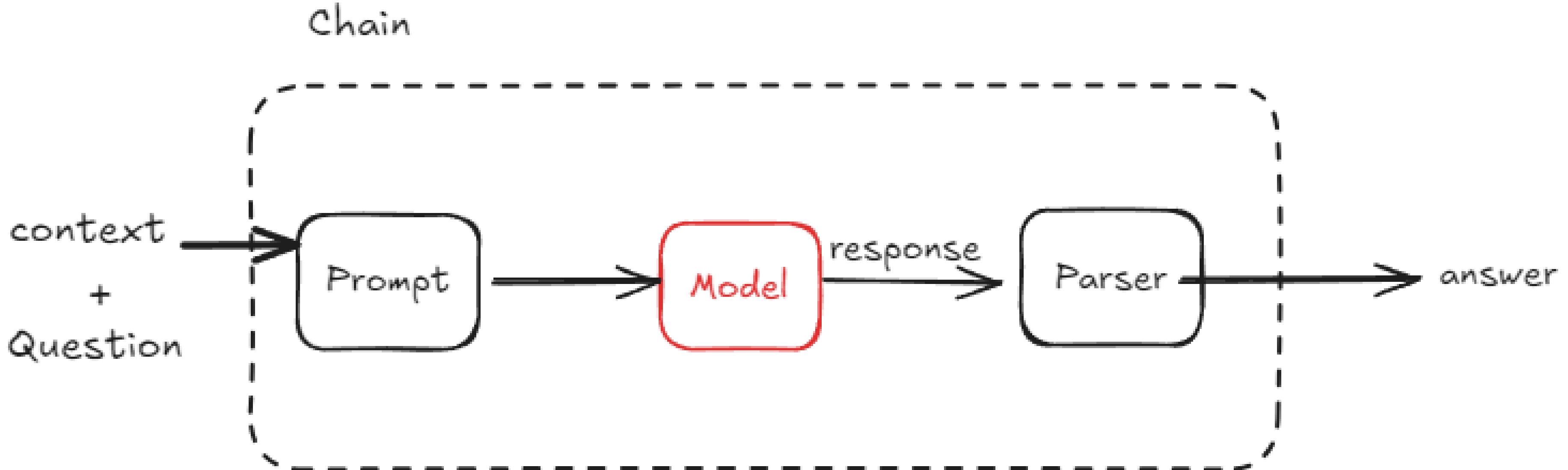
QUANTIZATION EXAMPLE

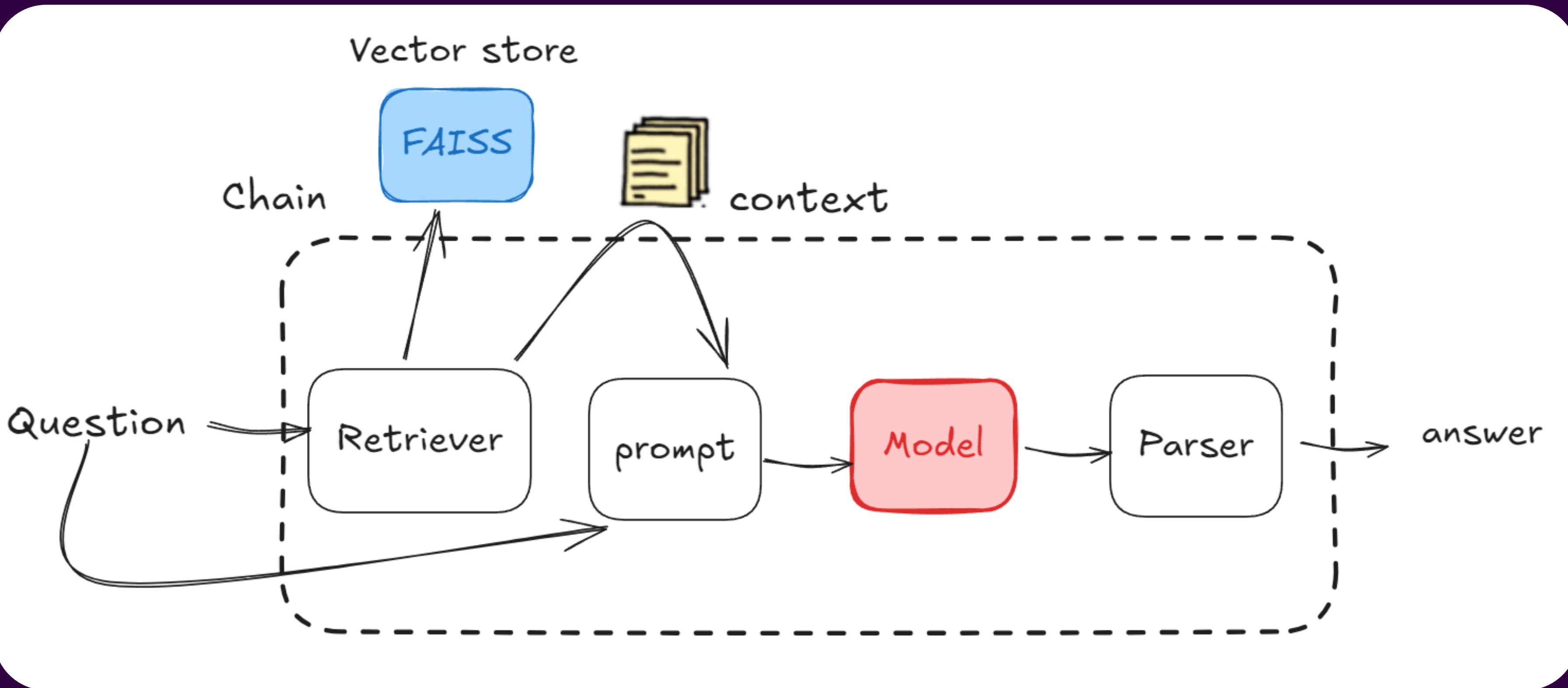
```
client.recreate_collection(  
    collection_name=collection_name,  
    vectors_config=models.VectorParams(  
        size=512, # dimension of your vectors #PQ works best with float32 vectors and dimensions divisible by 4 or 8.  
        distance=models.Distance.COSINE,  
        quantization_config=models.QuantizationConfig(  
            quantization=models.Quantization(  
                quantization=models.ProductQuantization( # using PQ method  
                    enabled=True,  
                    compression="x8", # or "x4", "x16" # x8 means 8x compression  
                    always_ram=False, # keep quantized vectors on disk  
                )  
            )  
        )  
    )  
)
```

RAG architecture









How to get started ?

LangChain

- LangChain is a framework designed to simplify the creation of applications using large language models.



Llamaindex



- Llamaindex is a handy tool that acts as a bridge between your custom data and large language models (LLMs) which are powerful models capable of understanding human-like text.



DEMO

→ Question: mention all the managers of AINS this year
Answer is : The managers of AINS mentioned are:

1. Fayed Zouari - Project Manager
2. Makki Aloulou - Program Manager
3. Kacem Mathlouthi - Technical Manager
4. Aziz Amari - manages the Poster Session (not explicitly titled as a manager, but has a management role)
5. Ahmed Amin Chabbah - oversees community engagement (not explicitly titled as a manager, but has a management role)
6. Rayen Khammar - handles event content (not explicitly titled as a manager, but has a management role)

```
[ ] query= "who is the contact manager of AINS this year "
response = qa.run(query)
```

→ Question: who is the contact manager of AINS this year
Answer is : I don't know who the contact manager of AINS is this year. The provided context mentions various roles and team members, including Project

```
▶ query = "who is the ambassador from IEEE ESSTHS SB in this edition?"
response = qa.run(query)
```

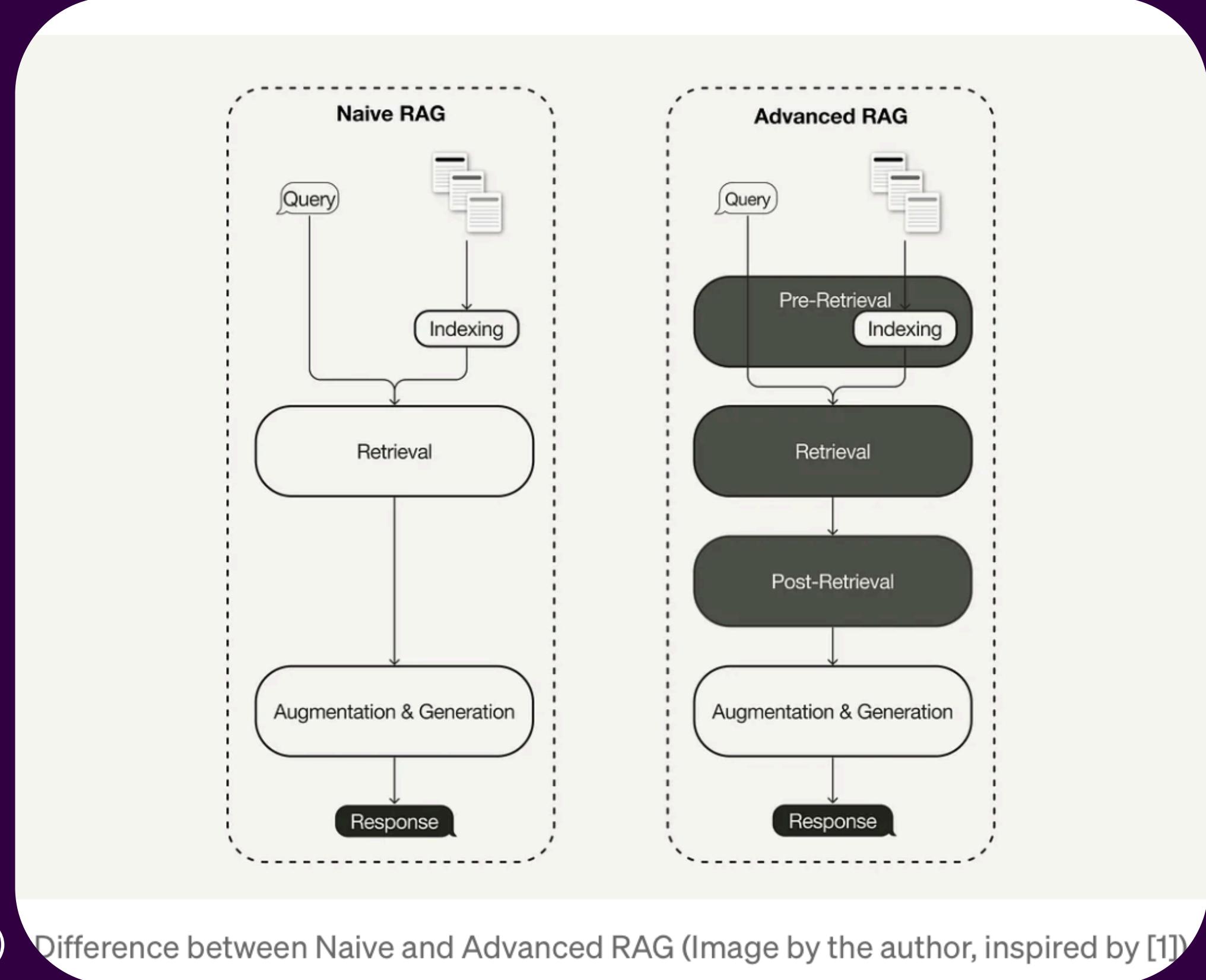
→ Question: who is the ambassador from IEEE ESSTHS SB in this edition?
Answer is : The ambassador from IEEE ESSTHS SB is Molka Attia.

```
[ ] query= "What talks will Hafedh Hichri and Mohamed Arbi Nsibi discuss in their talks in AINS 3.0 ?"
response = qa.run(query)
print(f"Question: {query}")
print(f"Answer is : {response}")
```

→ Question: What talks will Hafedh Hichri and Mohamed Arbi Nsibi discuss in their talks in AINS 3.0 ?
Answer is : Hafedh Hichri will discuss "Revolutionizing RAG with Chonkie", which is about his ultra-light Python library for efficient text chunking i

Naive RAG vs Advanced RAG

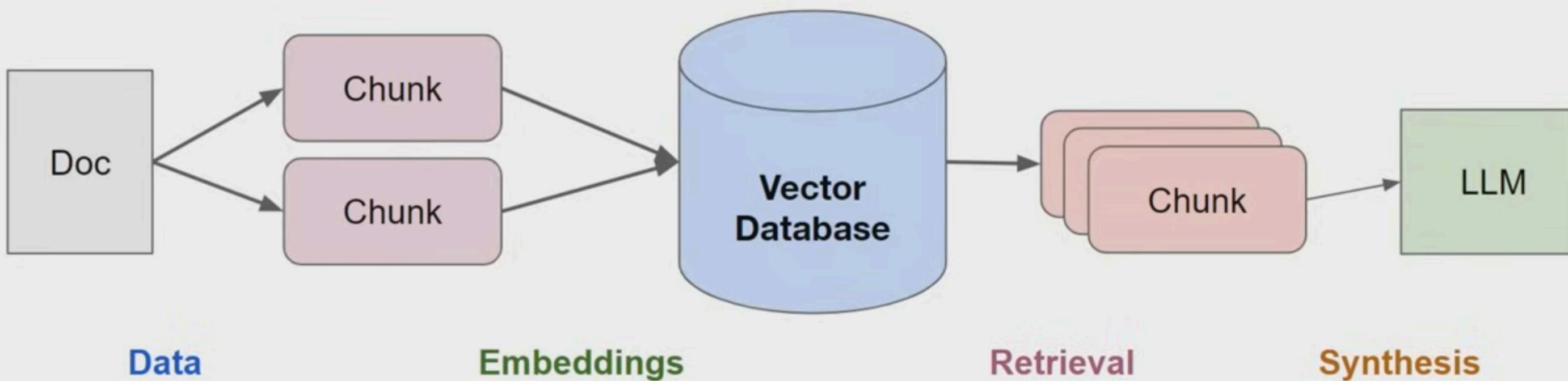
- There are many implementation to further improve performance of Naive RAG.
- Advanced RAG has evolved as a new paradigm with targeted enhancements to address some of the limitations of the naive RAG paradigm.
 - Advanced RAG techniques can be categorized into
 - pre-retrieval optimization,
 - retrieval optimization, and
 - post-retrieval optimization
- some examples :
 - Feedback loops (re-ranking, similarity score thresholds)
 - Hybrid Search (dense + keyword)
 - Contextual compression (summarize before feeding to LLM)
 - Multi-vector per chunk (dense embeddings per aspect)



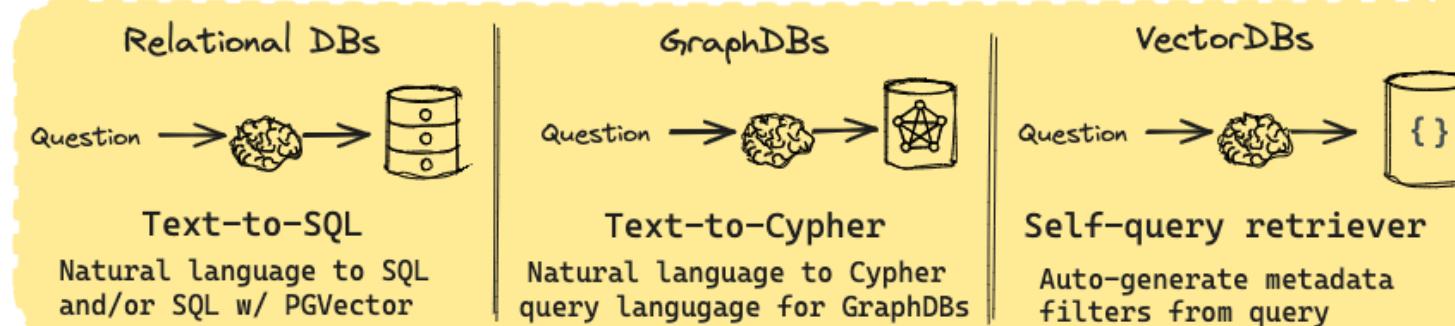
Naive RAG vs Advanced RAG

What do we do?

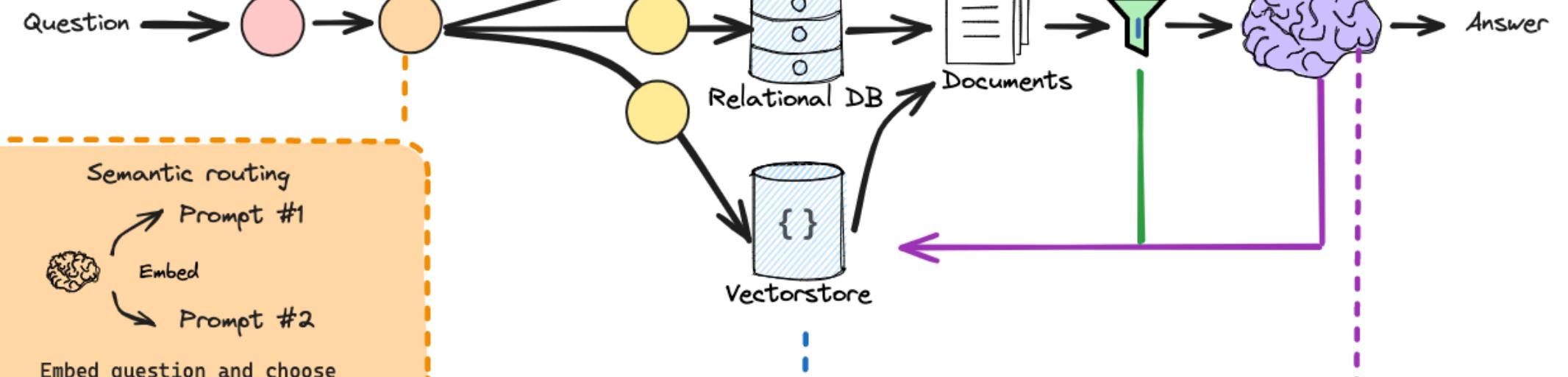
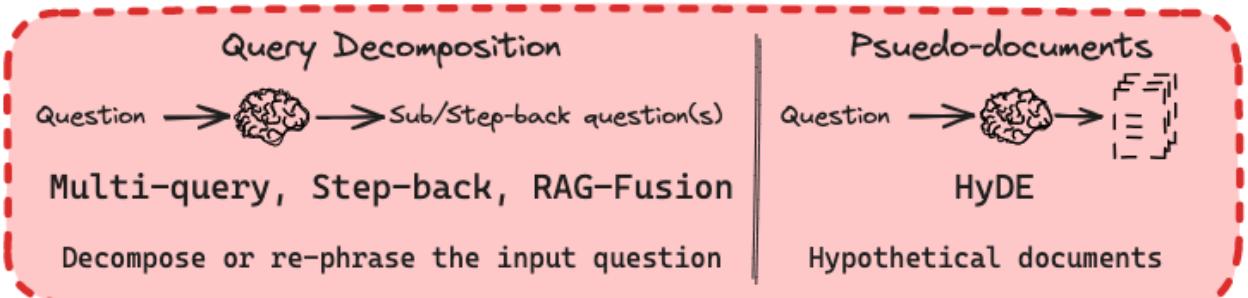
- **Data:** Can we store additional information beyond raw text chunks?
- **Embeddings:** Can we optimize our embedding representations?
- **Retrieval:** Can we do better than top-k embedding lookup?
- **Synthesis:** Can we use LLMs for more than generation? ✓



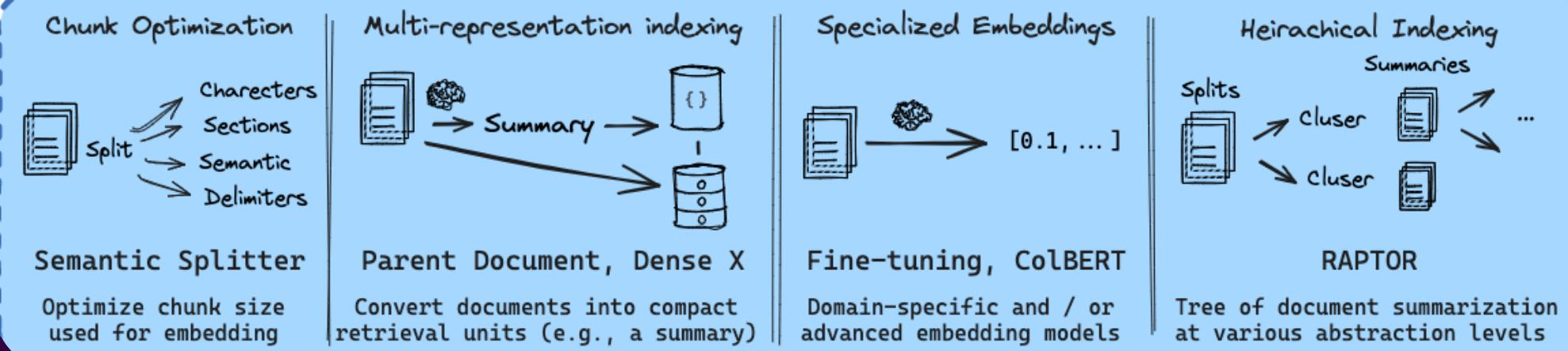
Query Construction



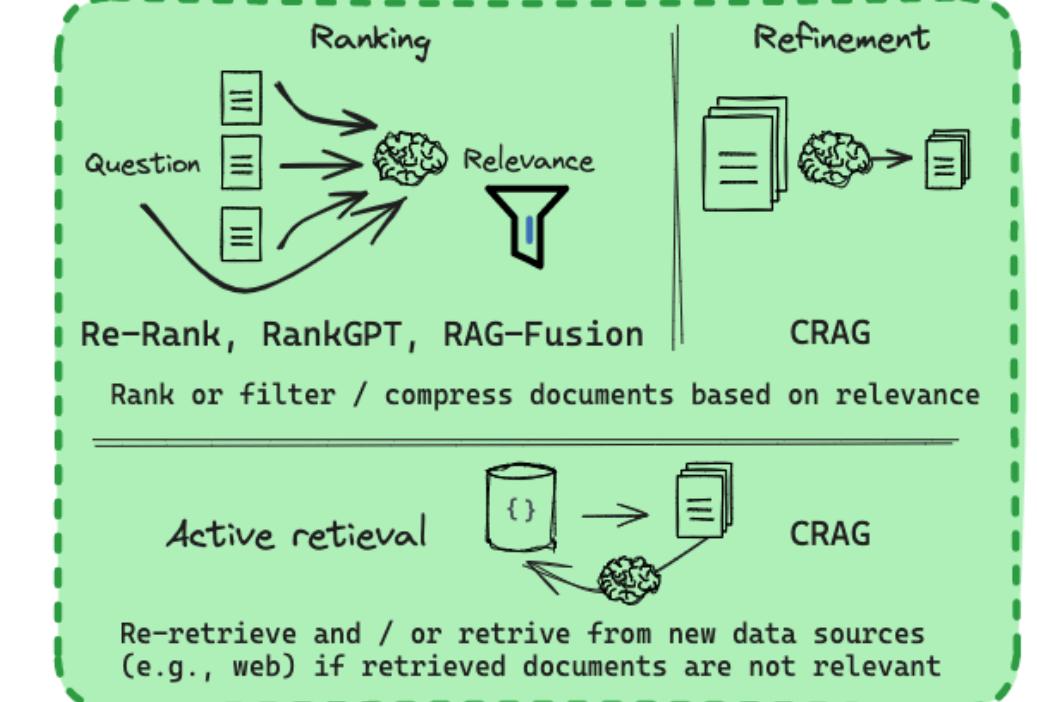
Query Translation



Indexing



Retrieval



Generation



Resources

- [All rag techniques](#)
- [Qdrant + DataTalks.club free course](#)
- [Similarity search HNSW](#)
- [Building nerual search service with ST and Qdrant](#)
- [Your RAG powered by Google Search Technology](#)
- <https://arxiv.org/abs/2005.11401>
- [Embedding models leaderboard](#)
- [Let's talk about LlamaIndex and LangChain](#)
- [Retrieval-Augmented Generation \(RAG\) framework in Generative AI](#)
- [Retrieval-Augmented Generation \(RAG\): From Theory to LangChain Implementation](#)
- [Advanced Retrieval-Augmented Generation: From Theory to LlamaIndex](#)
- [Retrieval-augmented generation for large language models: A survey \[arXiv\]](#)

**THANK YOU FOR YOUR
ATTENTION!!**

Q&A



<https://www.linkedin.com/in/mohammed-arbi-nsibi-584a43241/>