

Boas vindas

DANIEL SÓRIA
in



Otimizando LLM's com RAG

DANIEL SÓRIA



O que são LLM's?

LLM's - Large Language Models

- Modelos de gerais treinados com uma vasta quantidade de dados. (Large)
- Aprender padrão linguísticos (Language)
- Realizar várias tarefas como tradução, classificação, geração de texto.



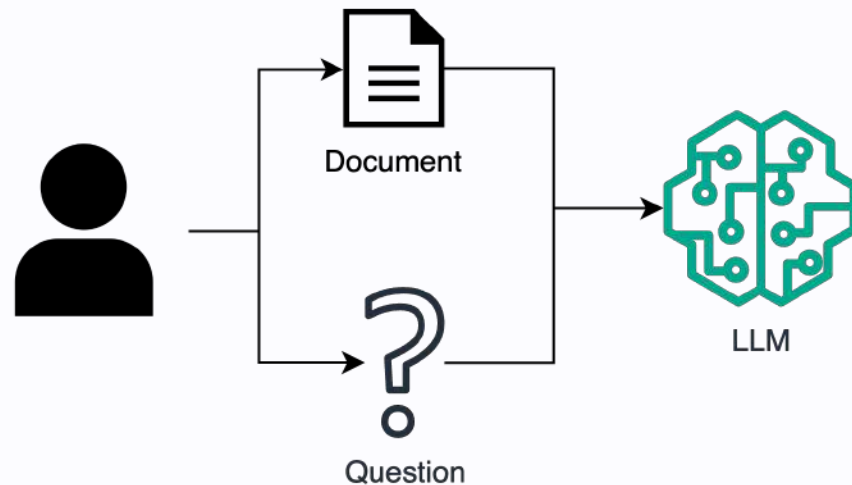
ChatGPT

LLaMA
by  **Meta**

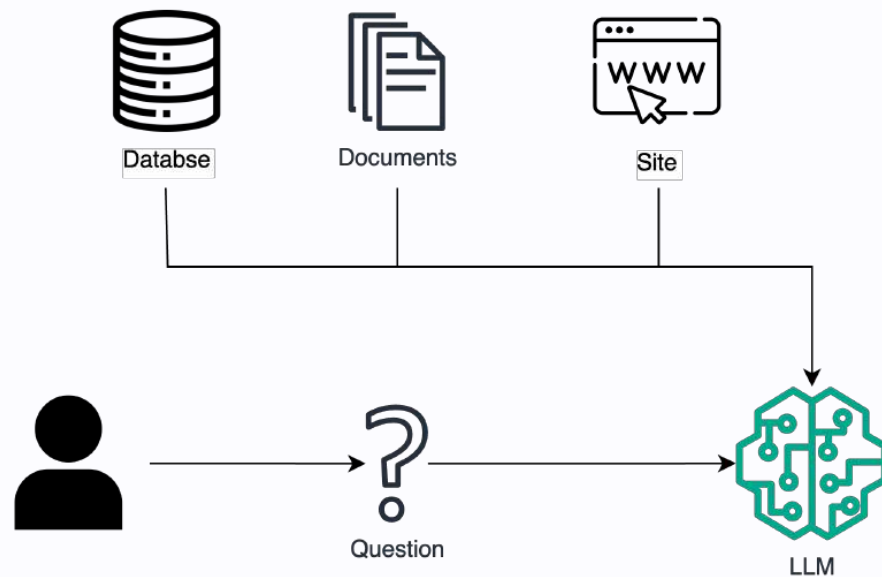
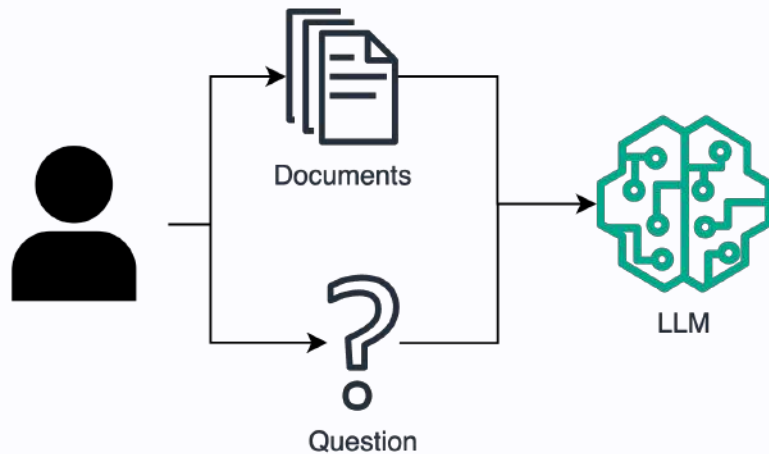

Gemini

LLM's - Large Language Models

- Limitados a base (gigantesca) de treinamento
- Dados públicos (ou não) gerais
- Não atendem necessidades específicas de cada empresa/ pessoa



LLM's - Large Language Models



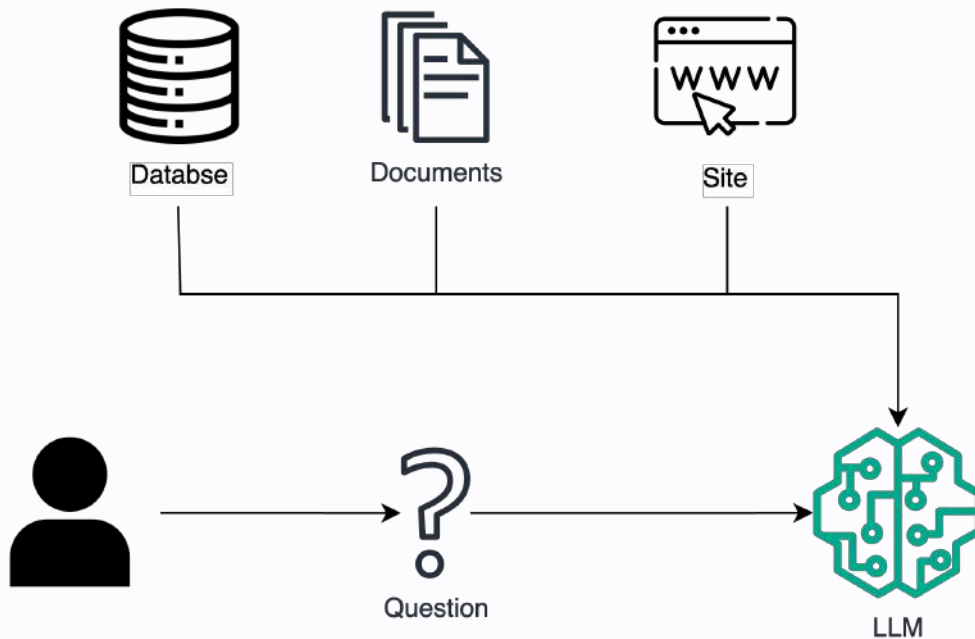
Otimizando LLM's com RAG

Retrieval

Augmented

Generation

Retrieval-Augmented Generation



Componentes do RAG

Retrieval-Augmented Generation

Componentes do RAG

Dados



Database



Document



Site



Question

Tecnologias



LLM



Embedding Model

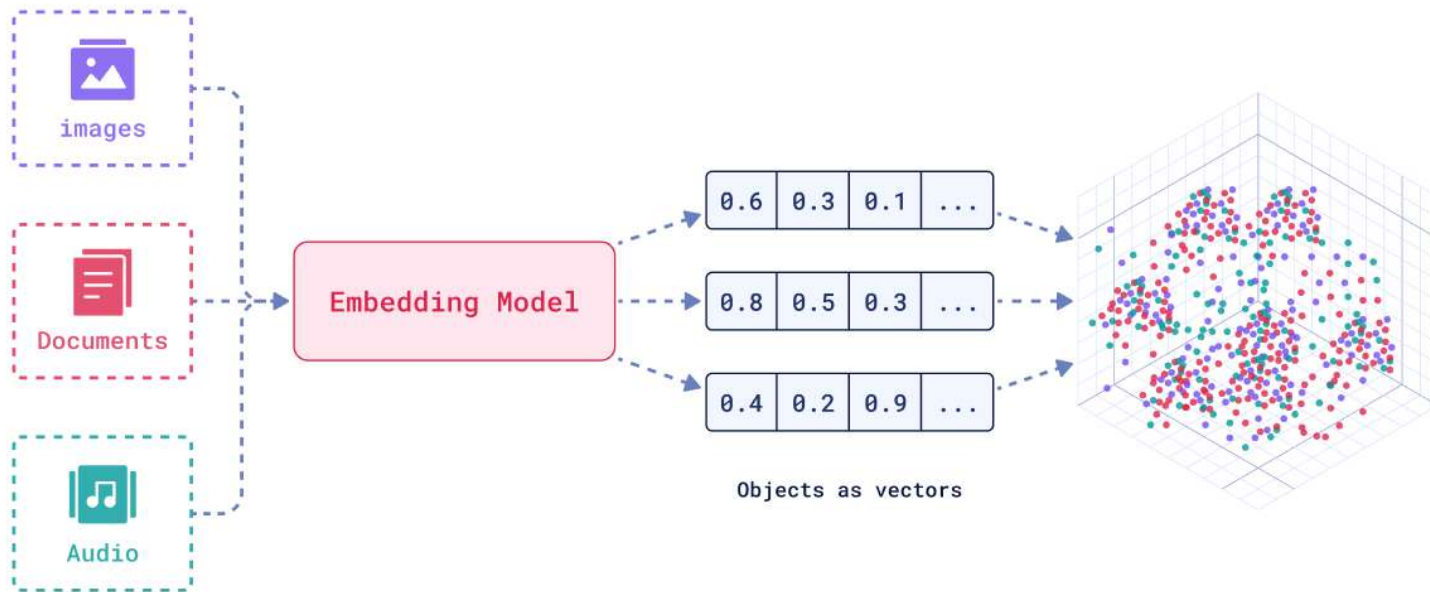


Vector DB

Componentes do RAG



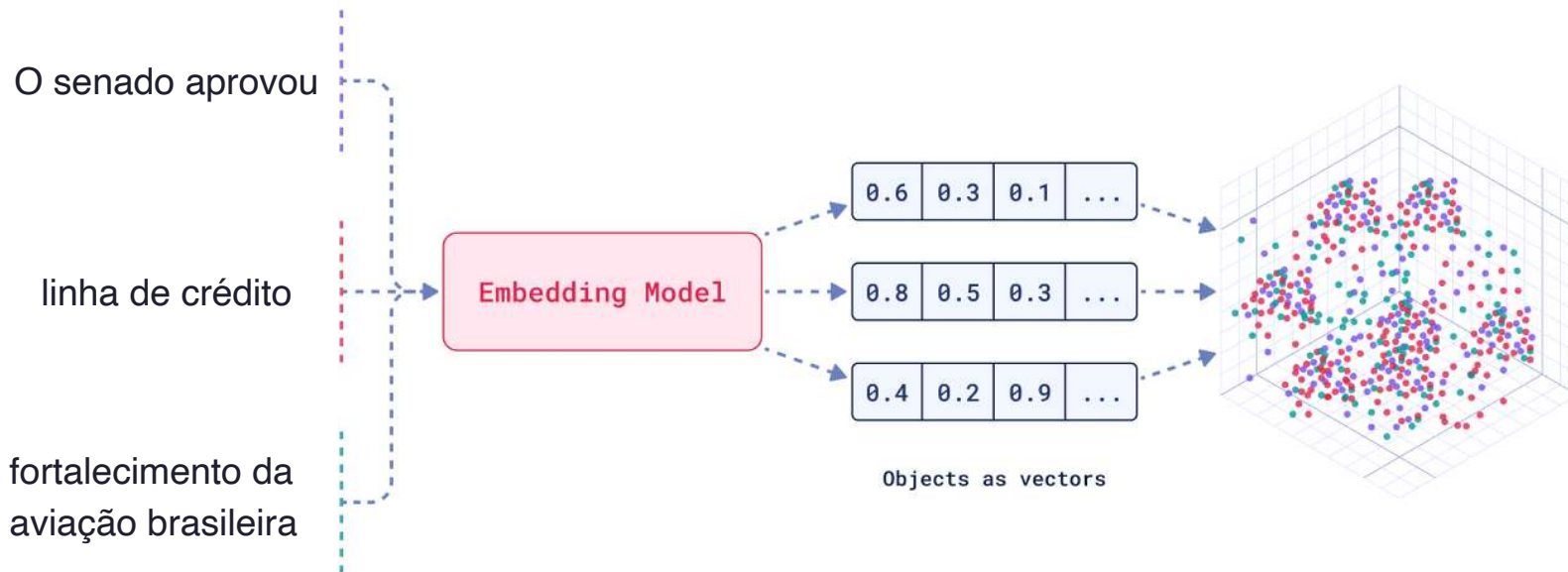
Embedding Model



Componentes do RAG



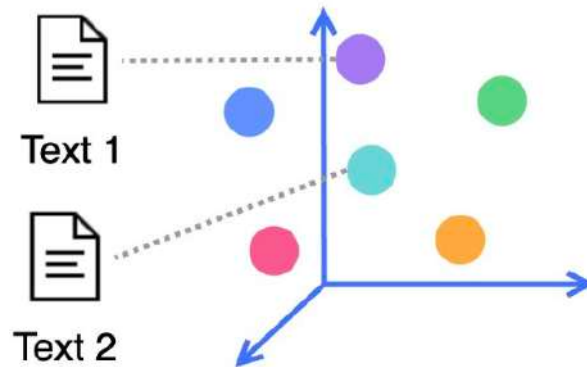
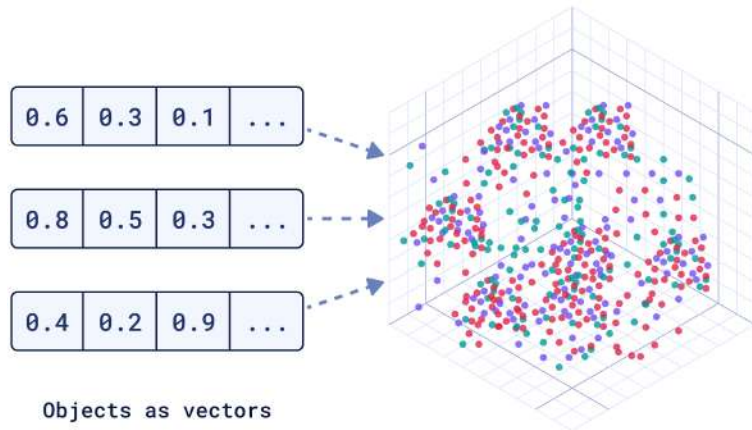
Embedding Model



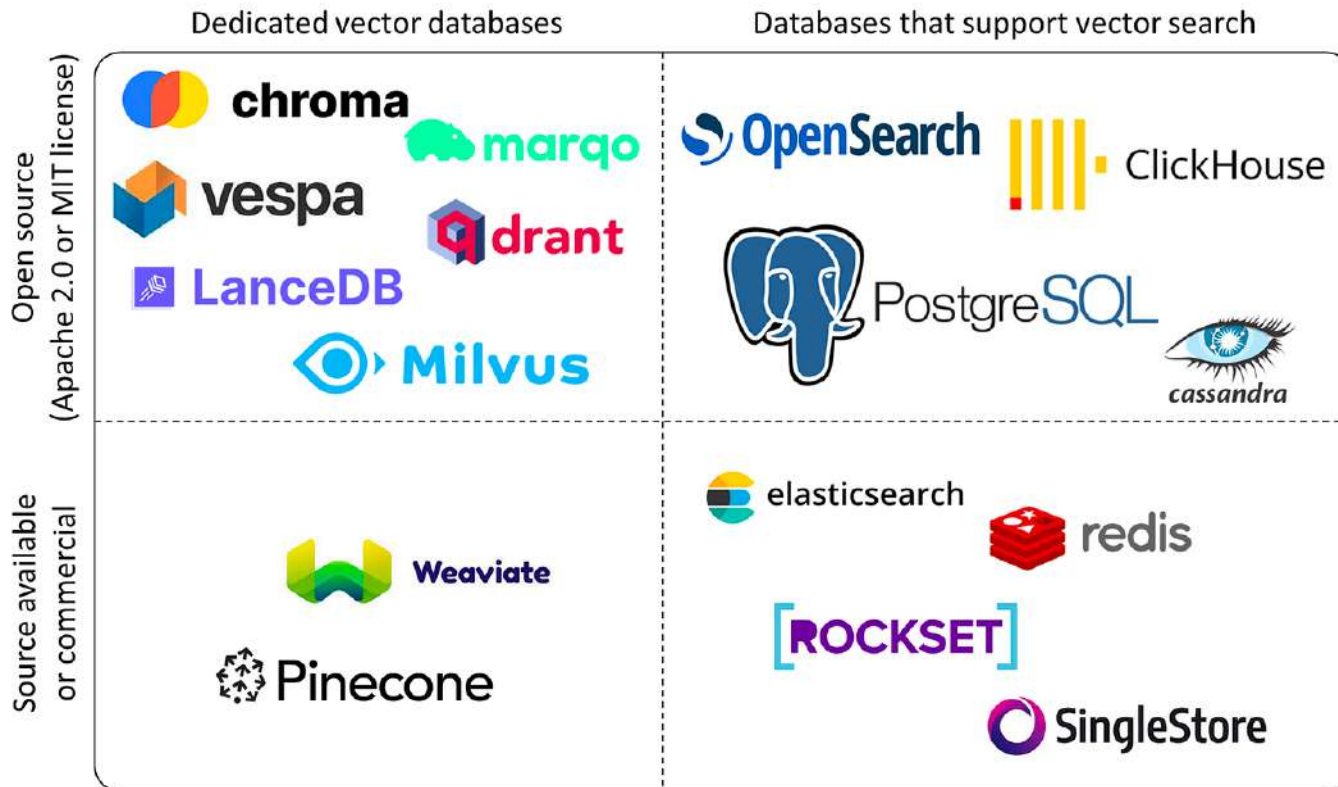
Componentes do RAG



Vector DB



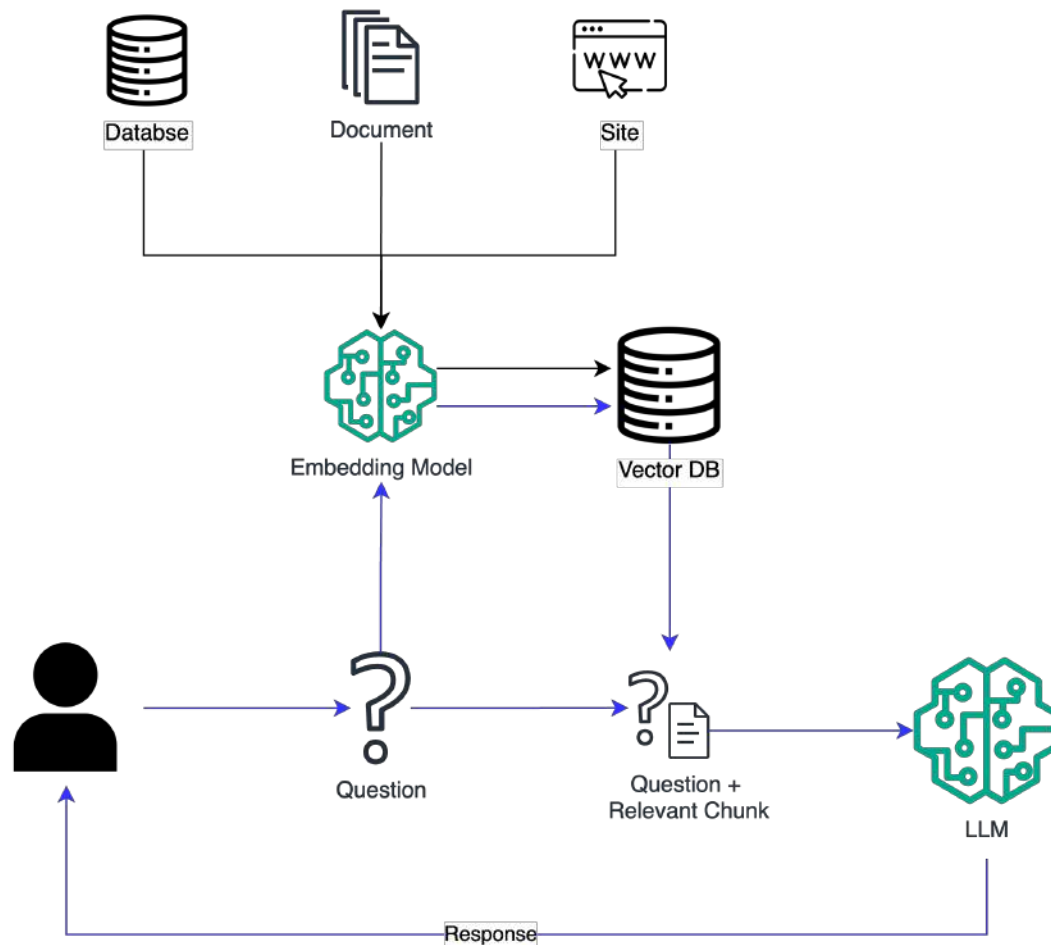
Componentes do RAG



Arquitetura de RAG

Retrieval-Augmented Generation

Arquitetura simples de RAG



Solução simples de RAG



Ask Your PDF

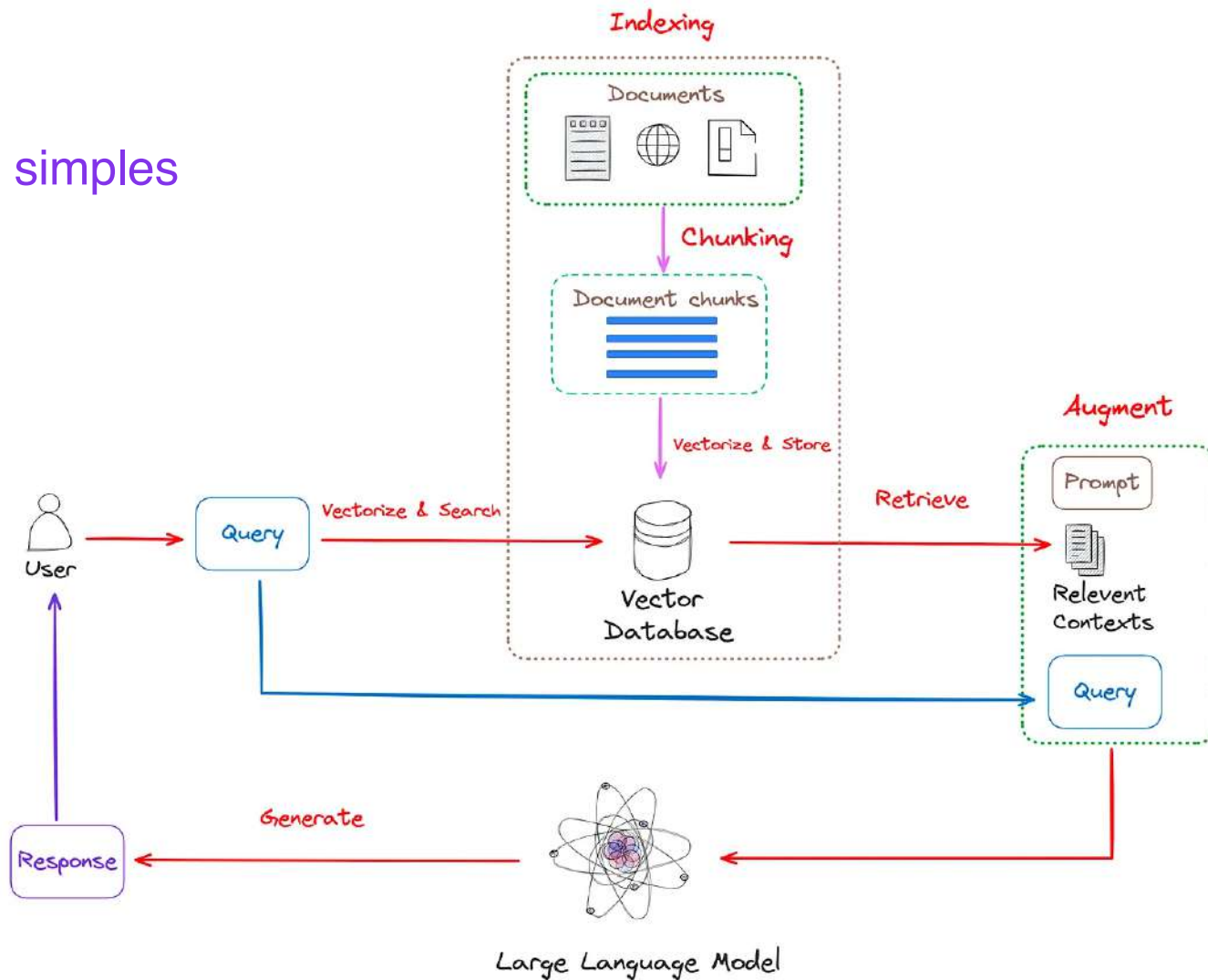


Chatbase

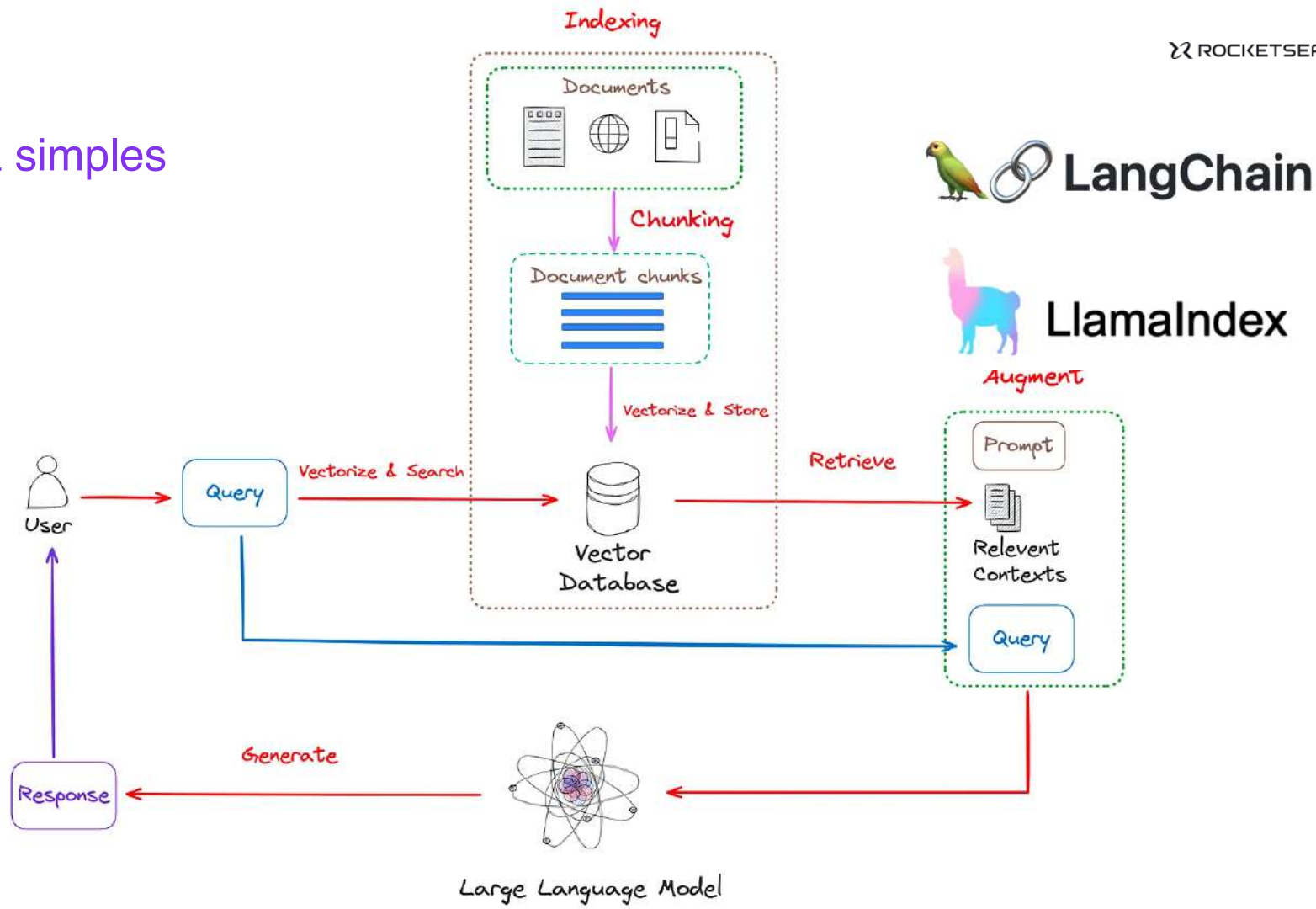


OpenAI

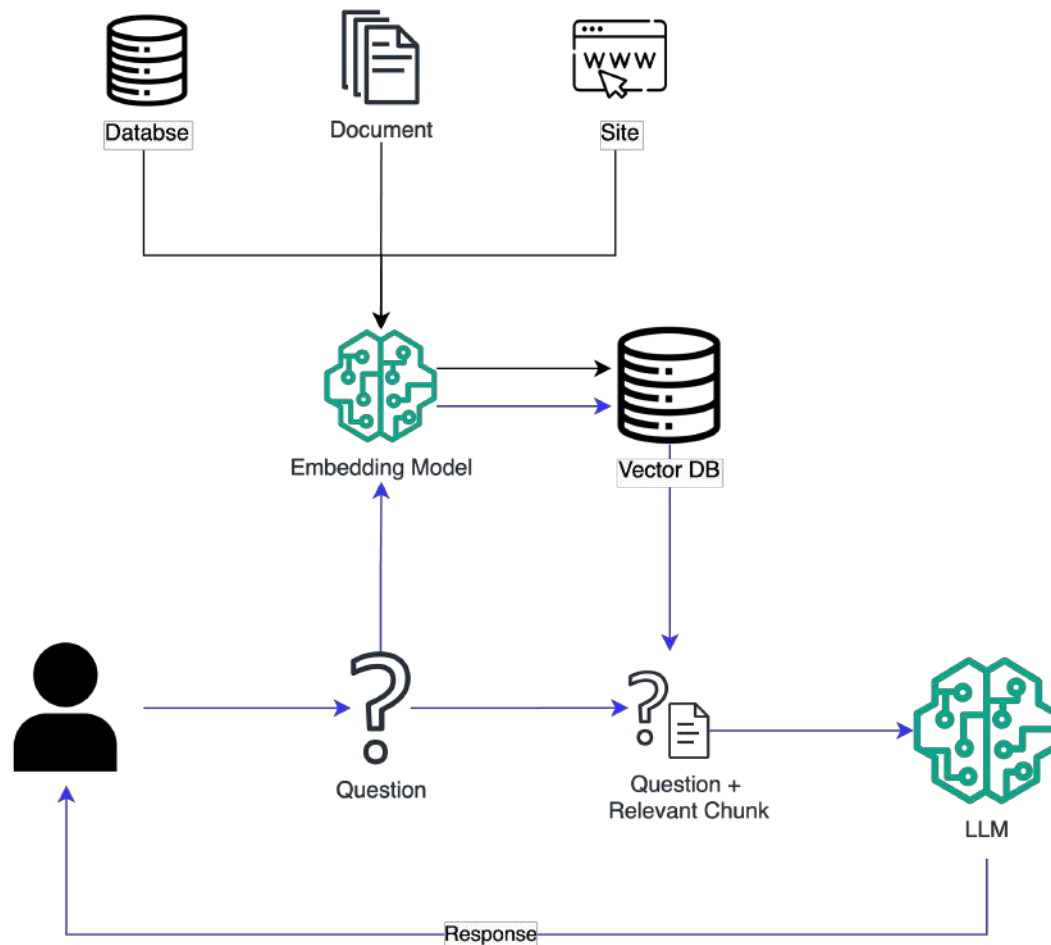
Arquitetura simples de RAG



Arquitetura simples de RAG

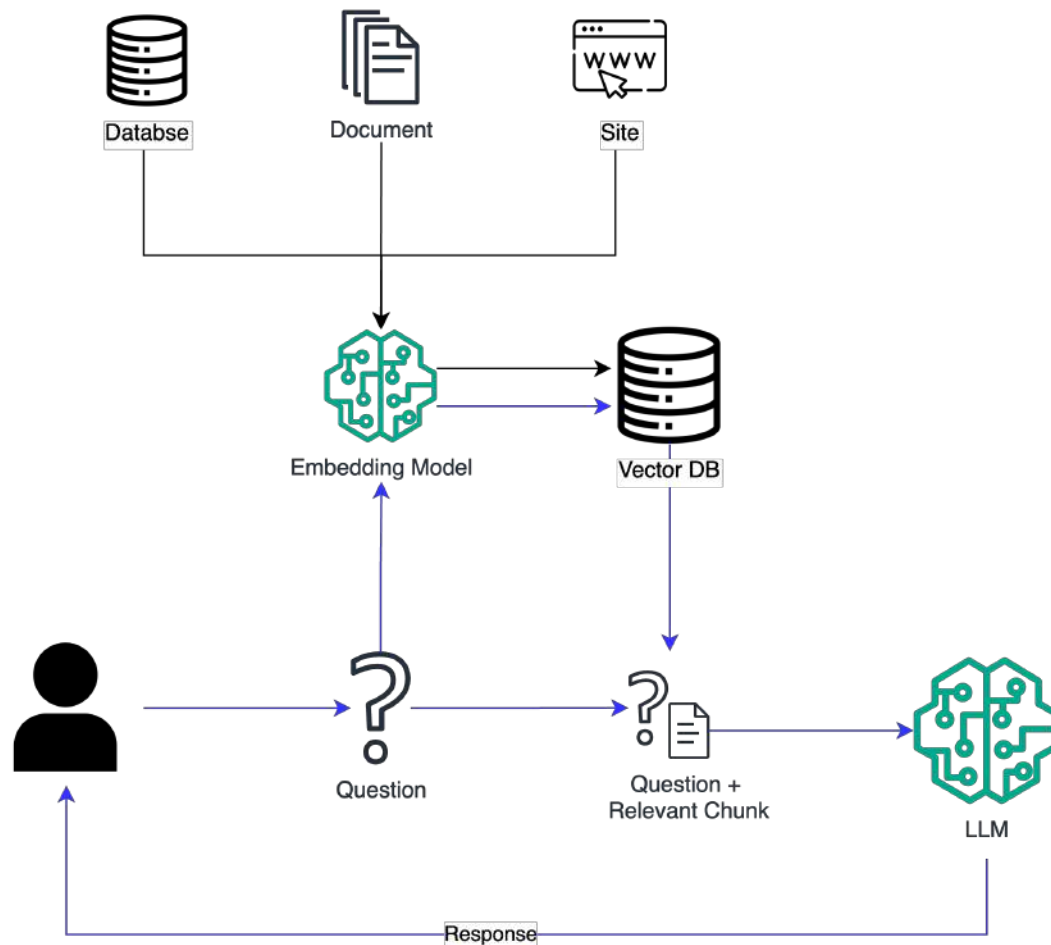


Arquitetura simples de RAG

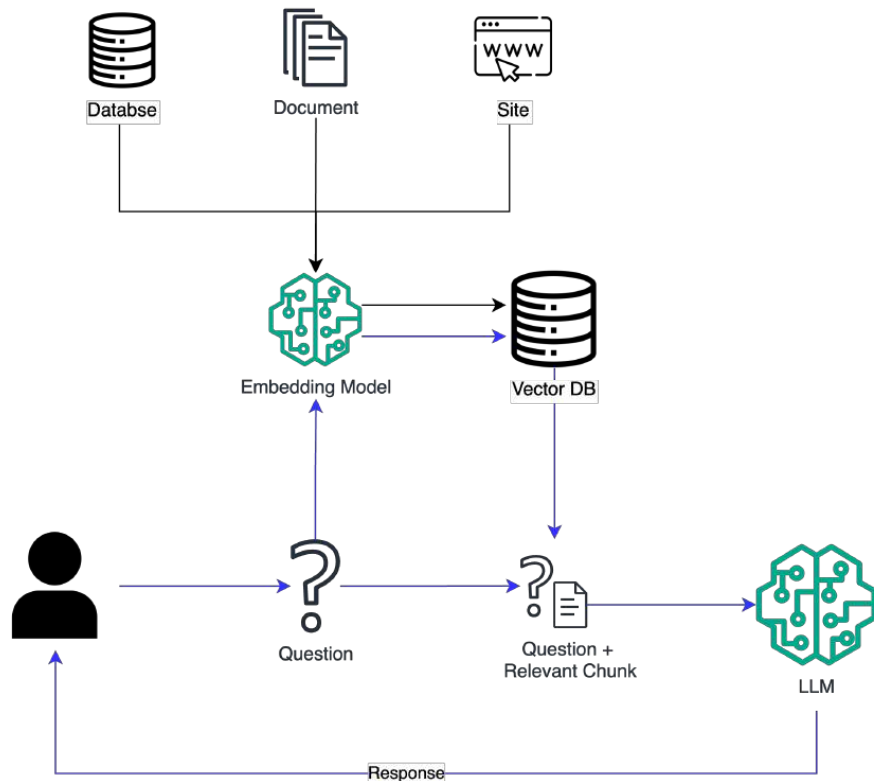


Desafios e limitações do RAG

Arquitetura simples de RAG



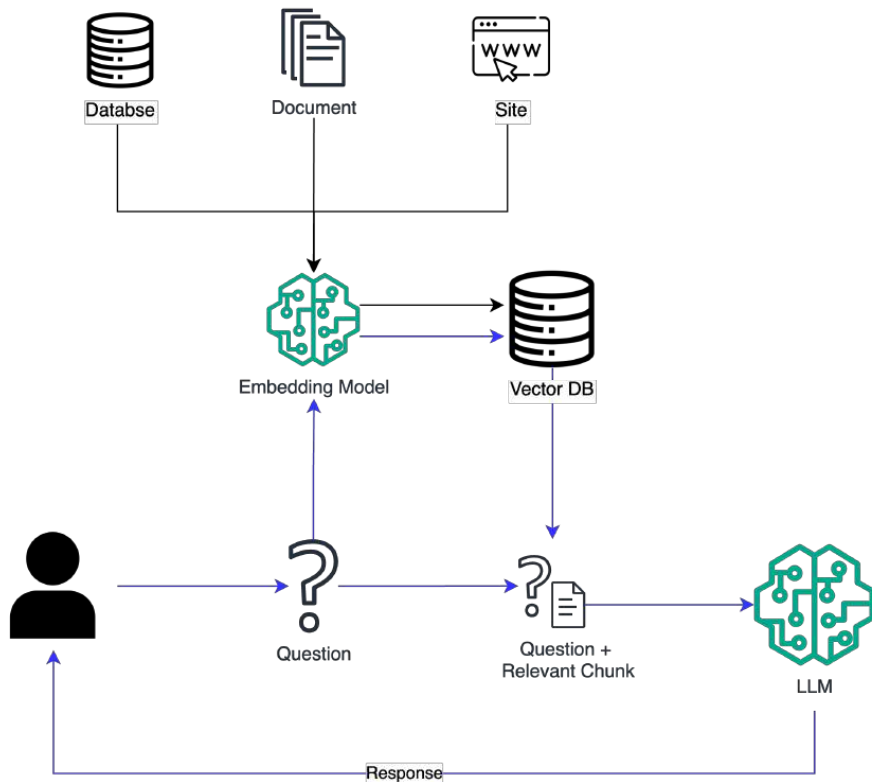
Arquitetura simples de RAG



1. Chunk size e Top-K;
2. Conhecimento de mundo;
3. Perda de informação;

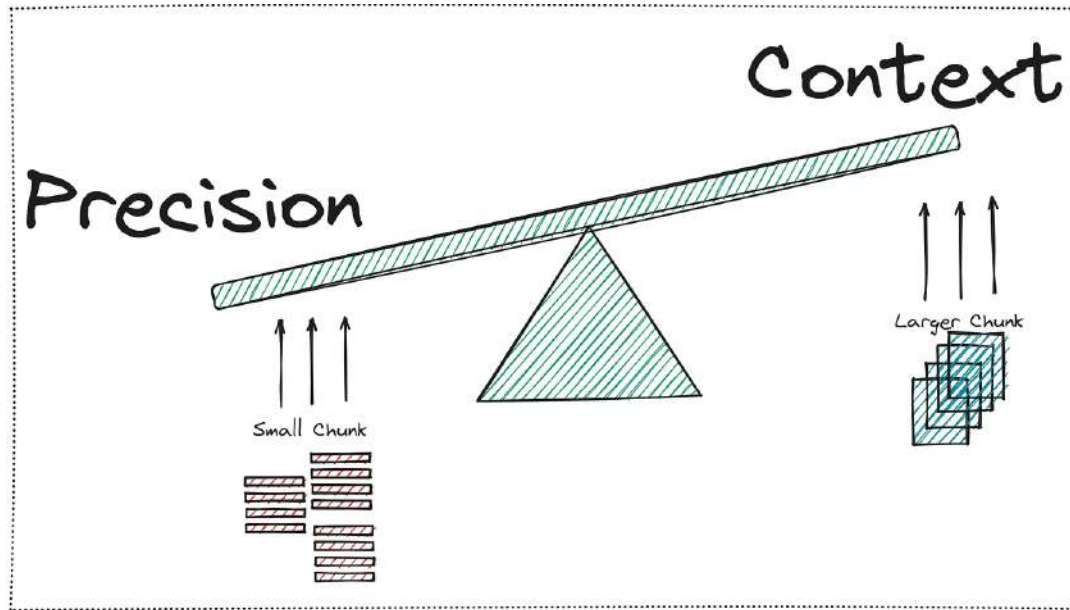
Advanced RAG

Naive RAG



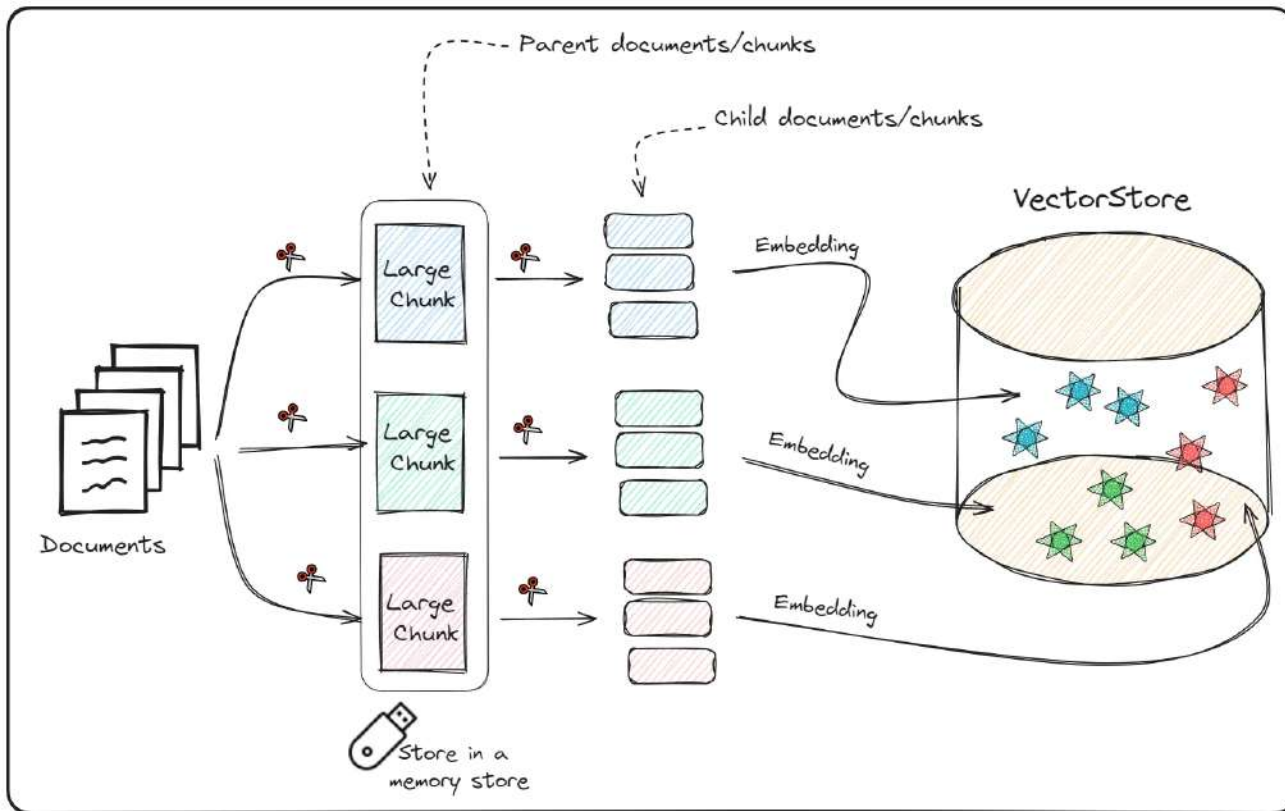
1. Chunk size e Top-K;
2. Conhecimento de mundo;
3. Perda de informação;

Naive RAG



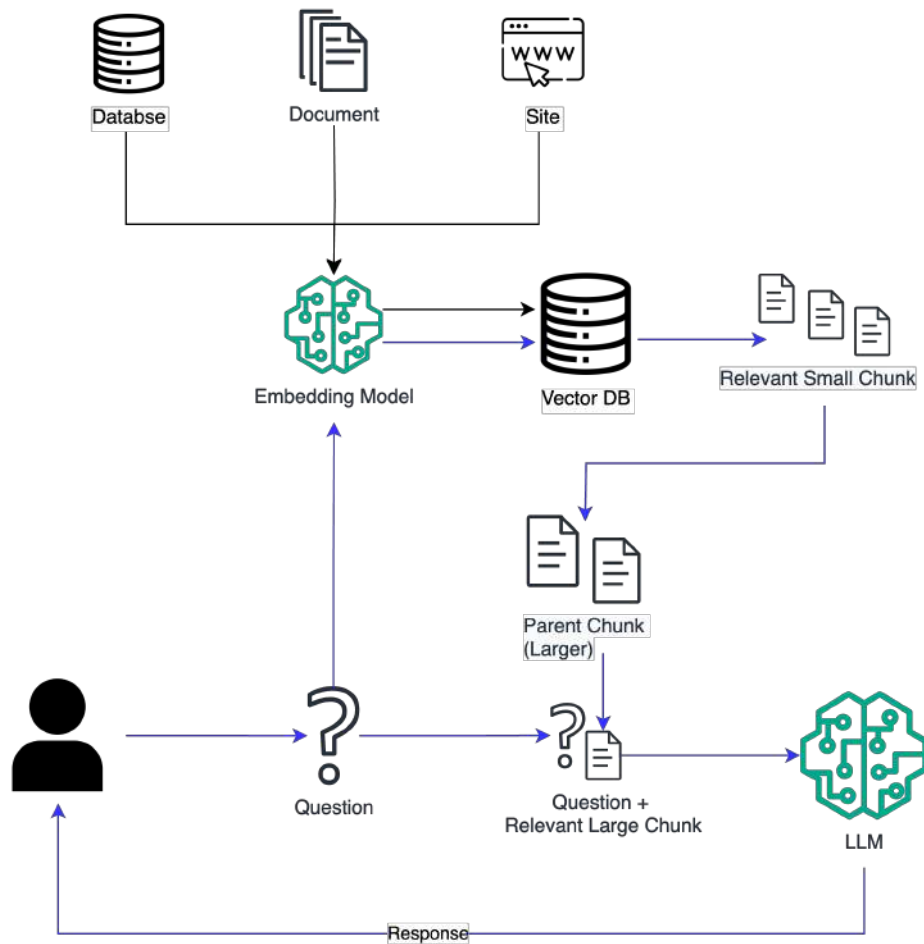
1. Chunk size e Top-K;
2. Conhecimento de mundo;
3. Perda de informação;

Parent Document Retriever

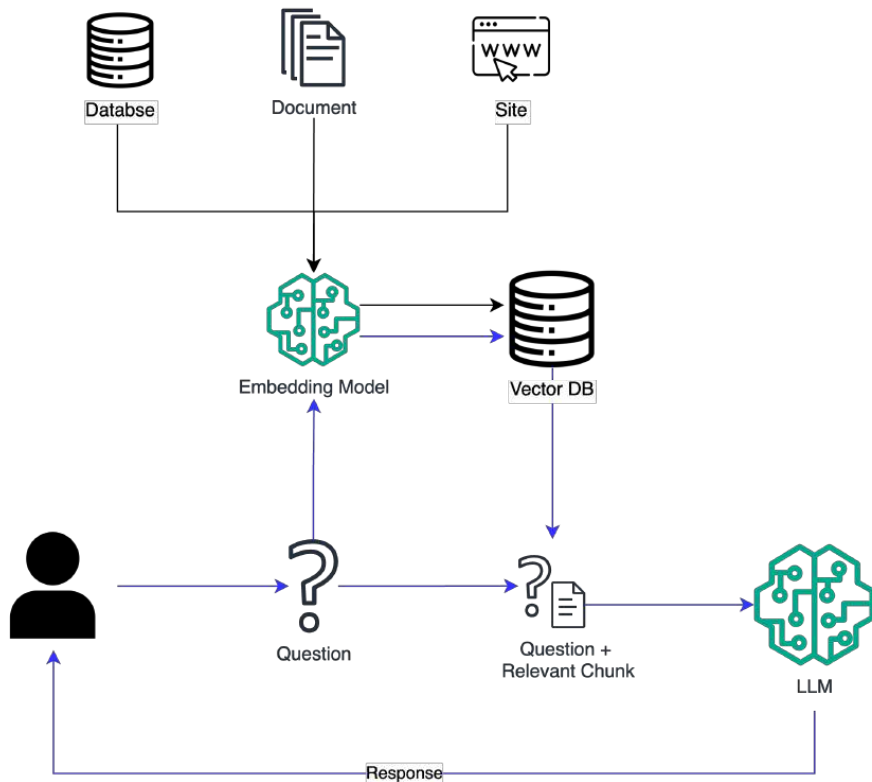


Parent Document RAG

Parent Document Retriever

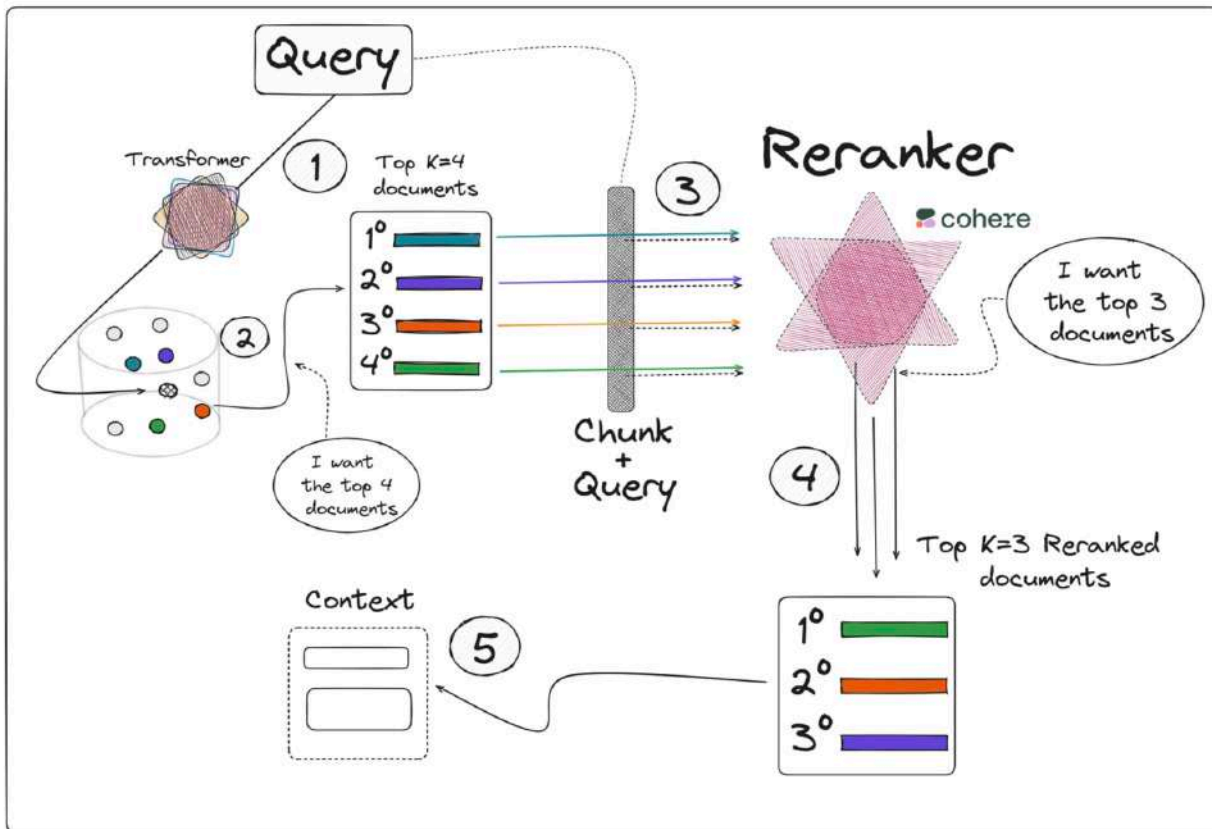


Naive RAG



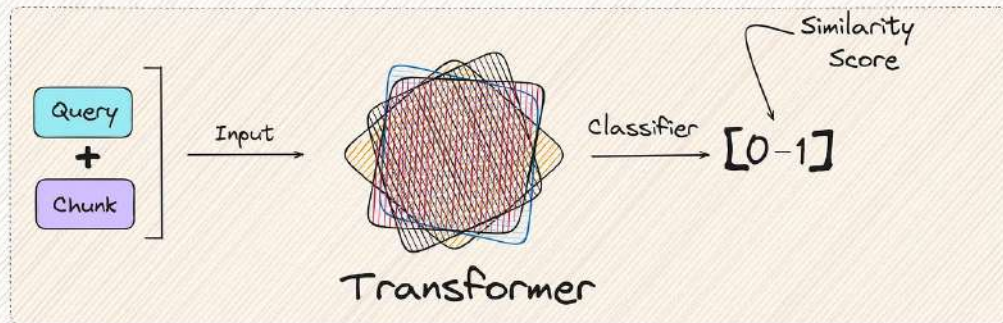
1. Chunk size e Top-K;
2. Conhecimento de mundo;
3. Perda de informação;

Contextual Compression RAG (Reranker)

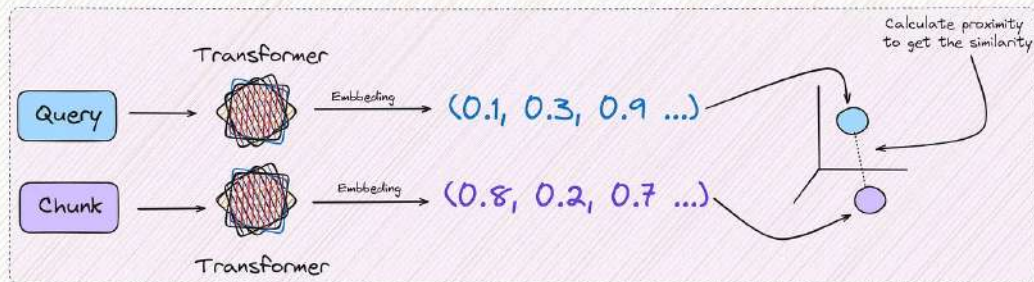


Contextual Compression RAG (Reranker)

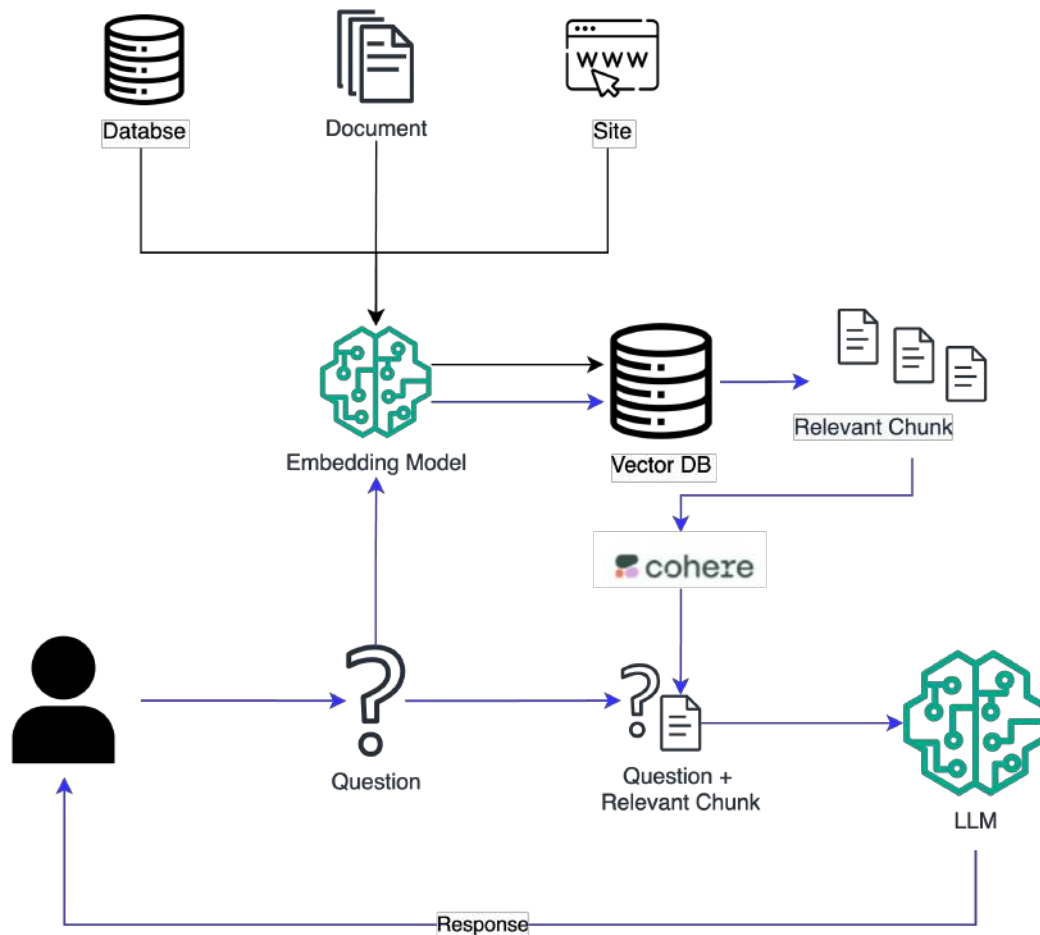
Cross-Encoder (It is use to the reranking)



Bi-Encoder (It is use to the naive retriever)

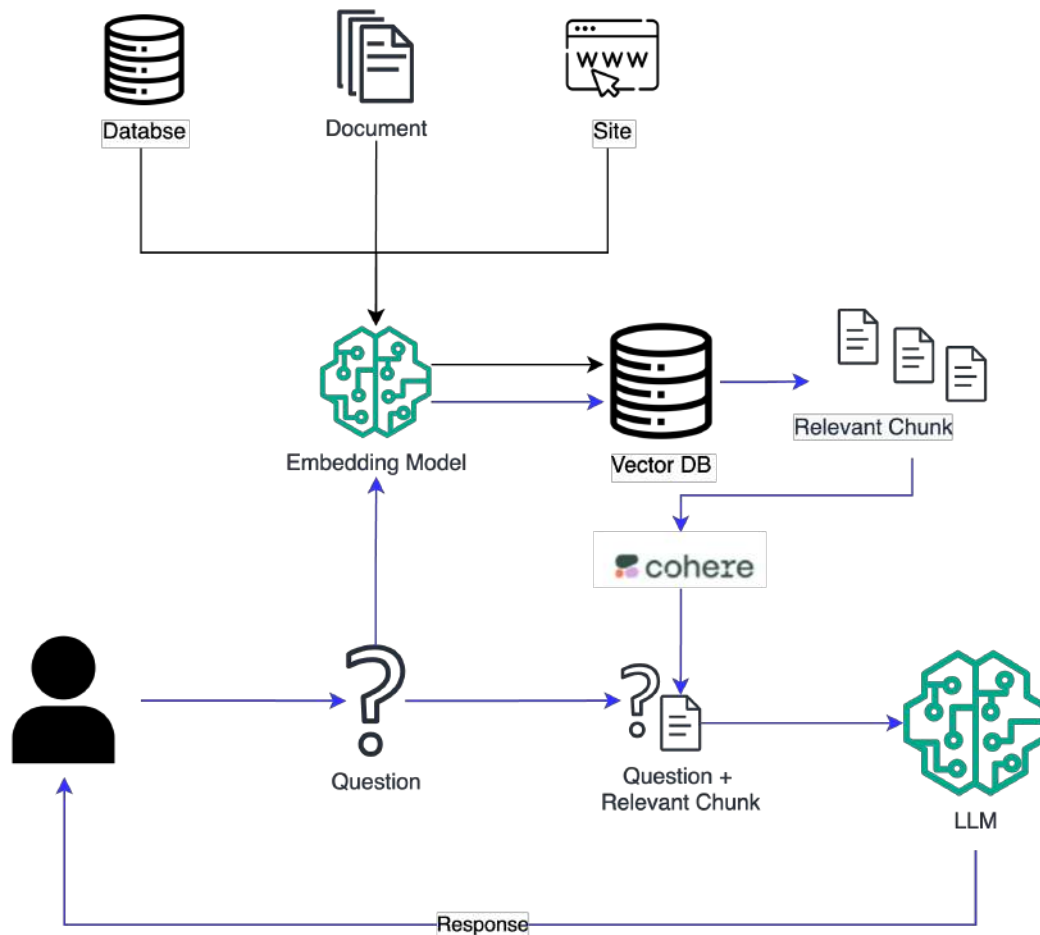


Contextual Compression RAG (Reranker)

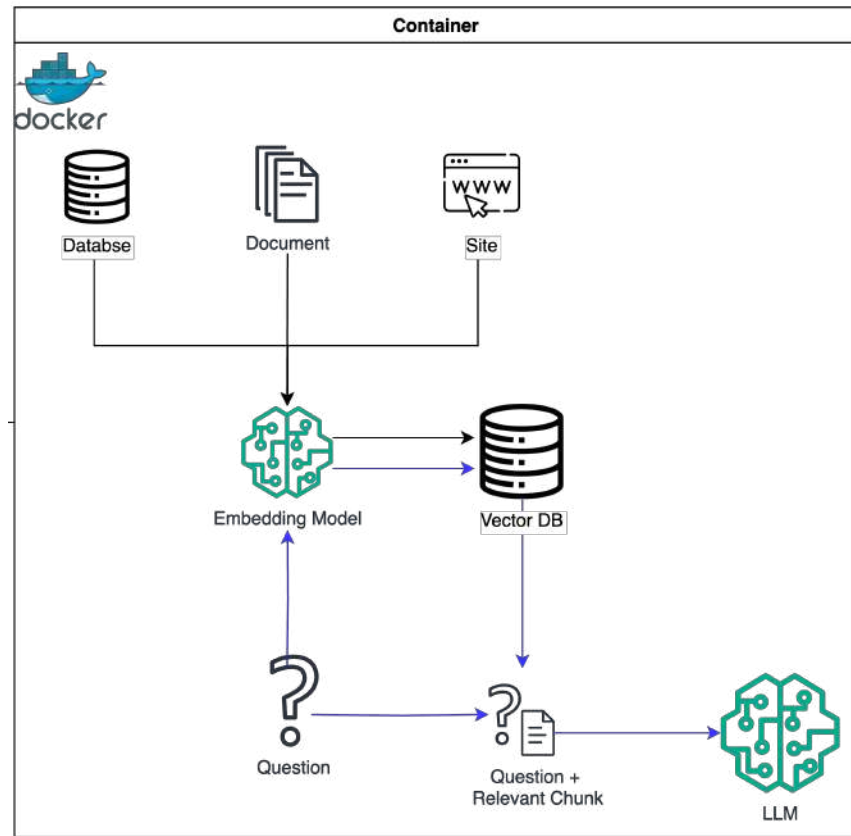


Deploy do RAG na Cloud

Arquitetura simples de RAG



Arquitetura simples de RAG



Arquitetura simples de RAG

