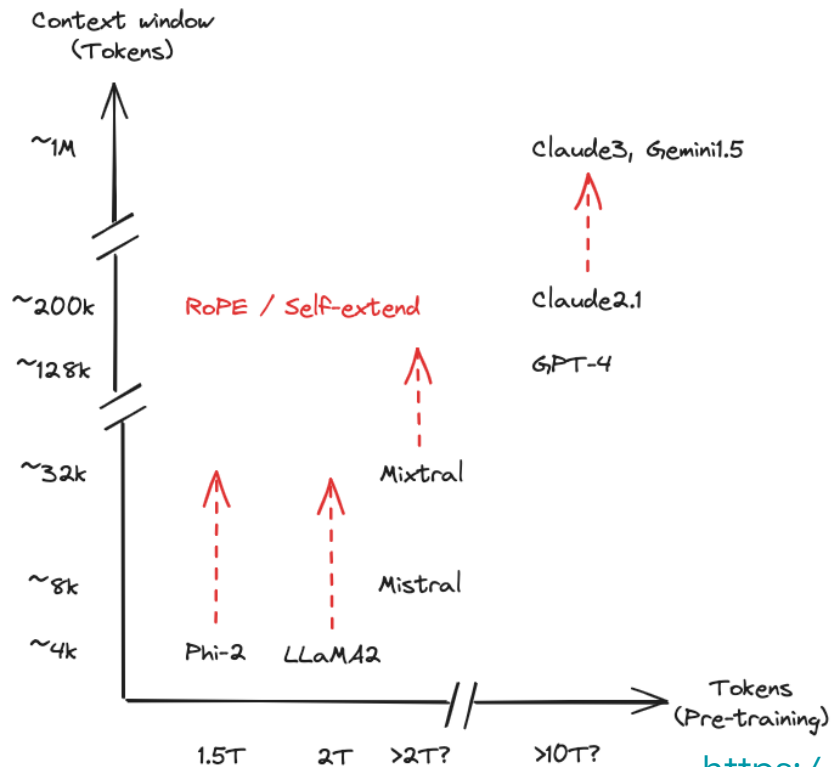


RAG from scratch: Overview

Lance Martin
Software Engineer, LangChain
[@RLanceMartin](https://twitter.com/RLanceMartin)

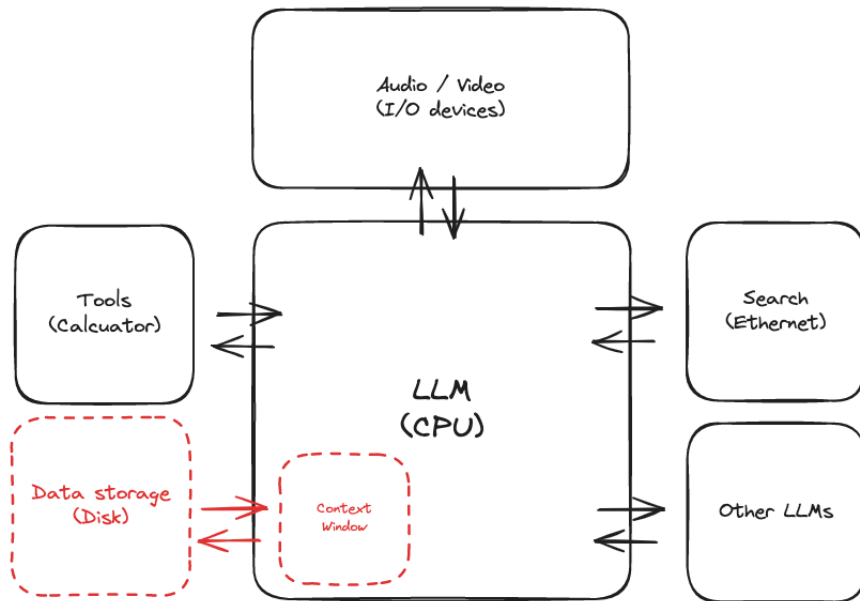
> 95% of the world's data is “private”, but we can “feed it” to LLMs



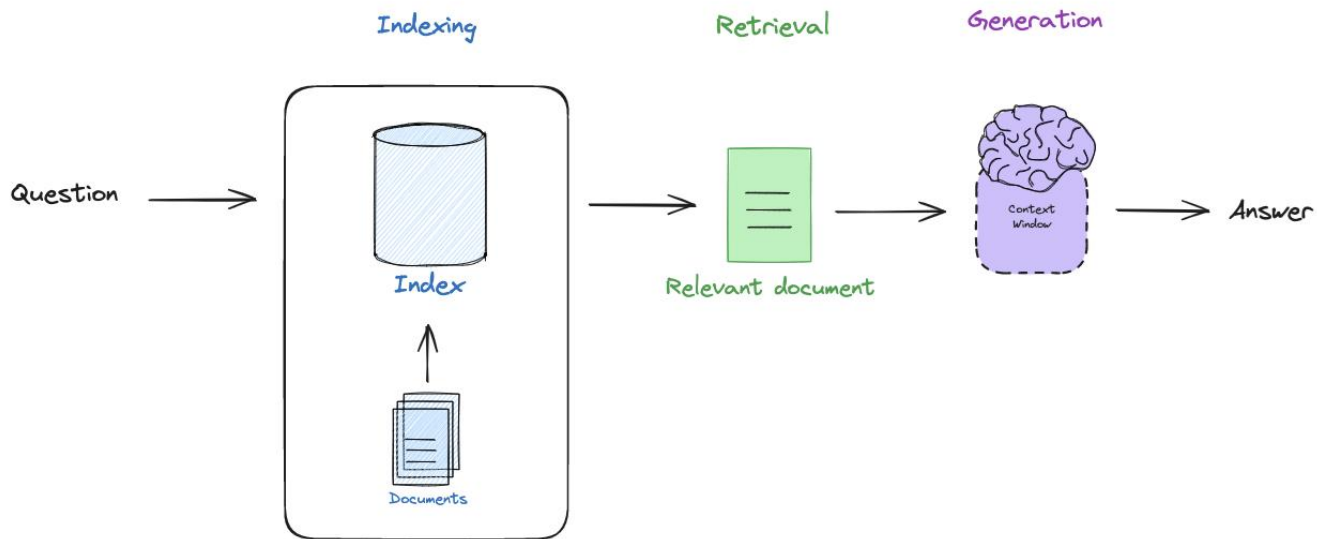
<https://huggingface.co/blog/mixtral>

<https://x.com/RihardJarc/status/1778082161595208124>

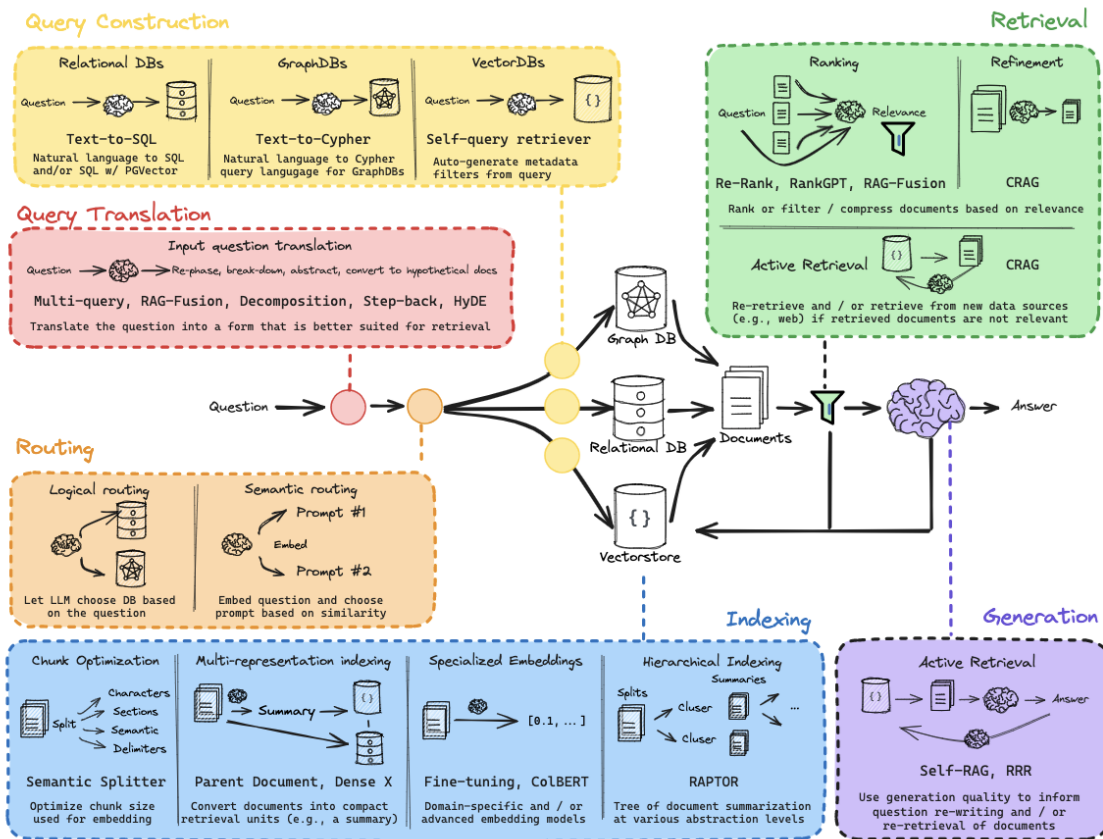
Connecting LLMs to external data is a central need



Retrieval Augmented Generation



We'll start from scratch, and build up to state of the art in RAG



Outline

Basics

- Indexing
- Retrieval
- Generation

Advanced

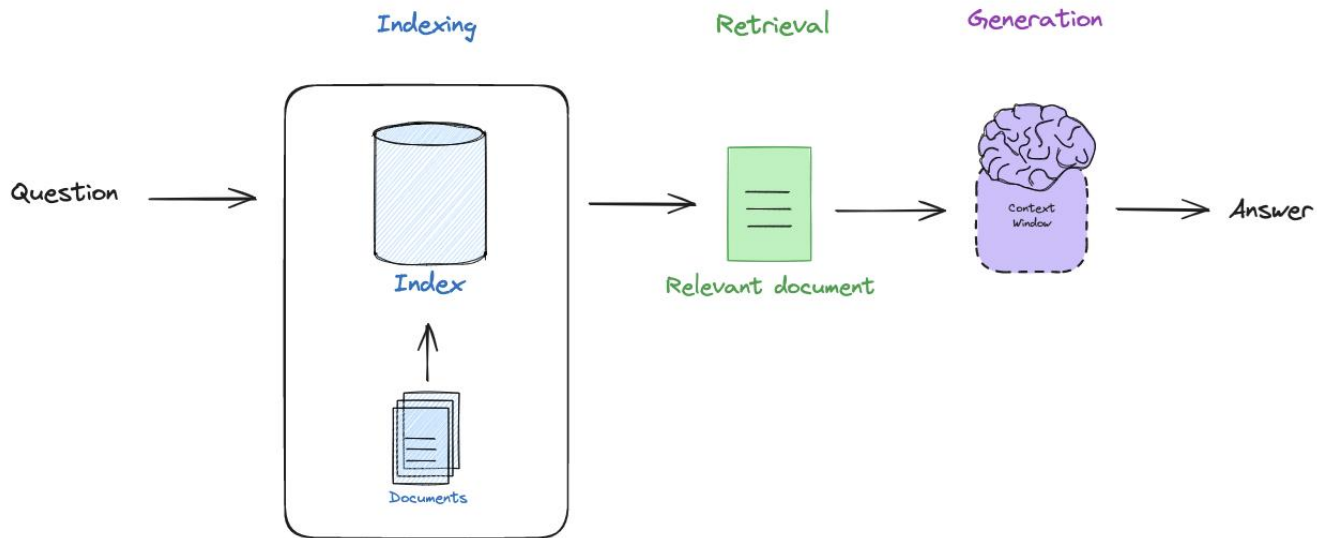
- Query transformations
- Routing
- Query construction
- Indexing
- Retrieval
- Generation

Code walk-through

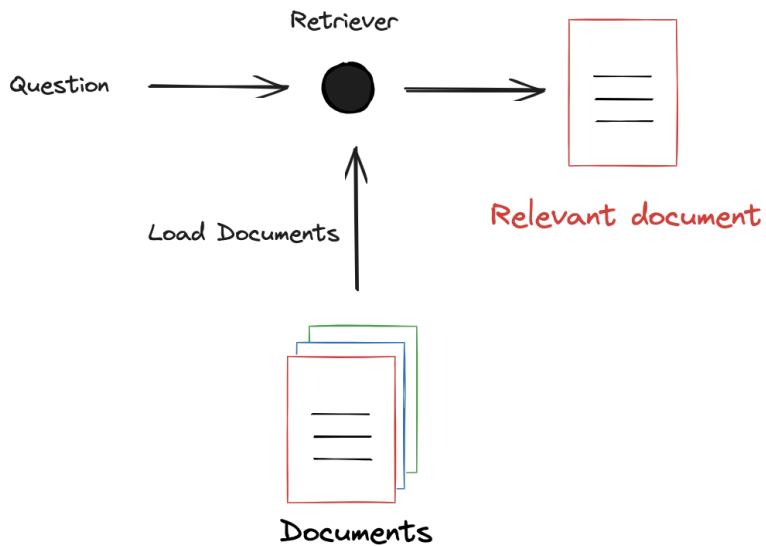
RAG from scratch: Indexing

Lance Martin
Software Engineer, LangChain
[@RLanceMartin](https://twitter.com/RLanceMartin)

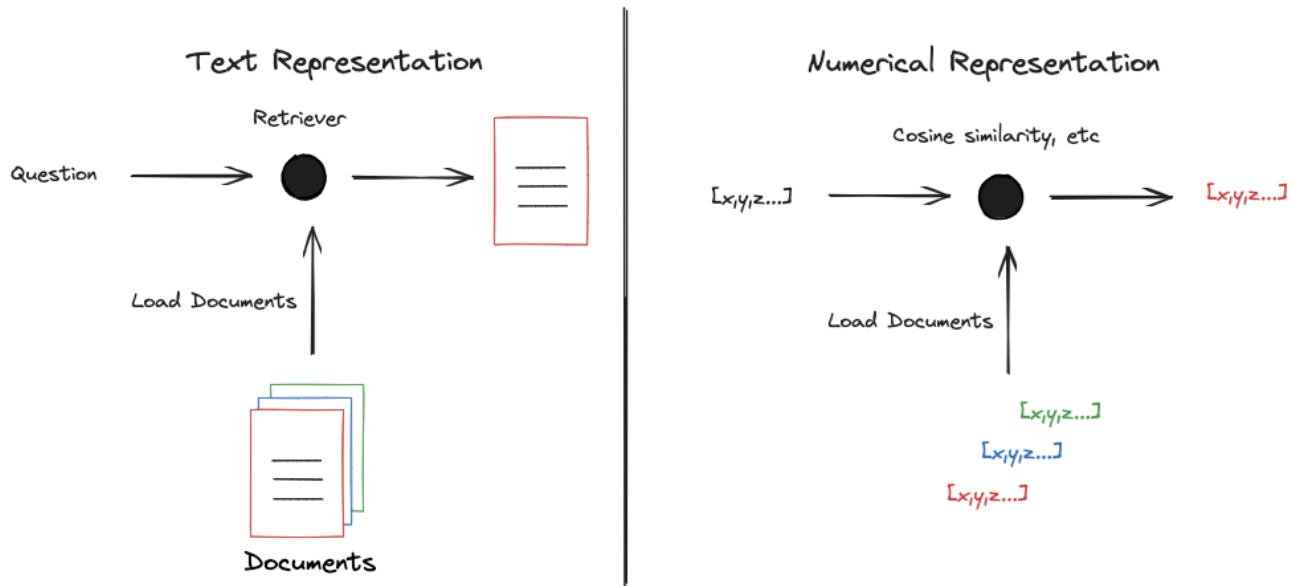
RAG motivation



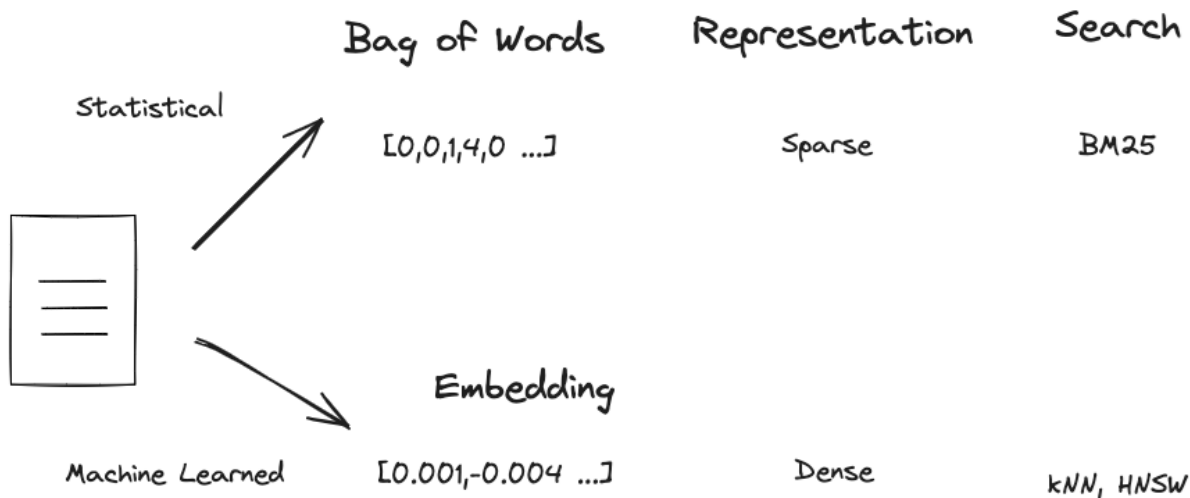
Document loading



Numerical representation for search



Statistical and machine learned representations

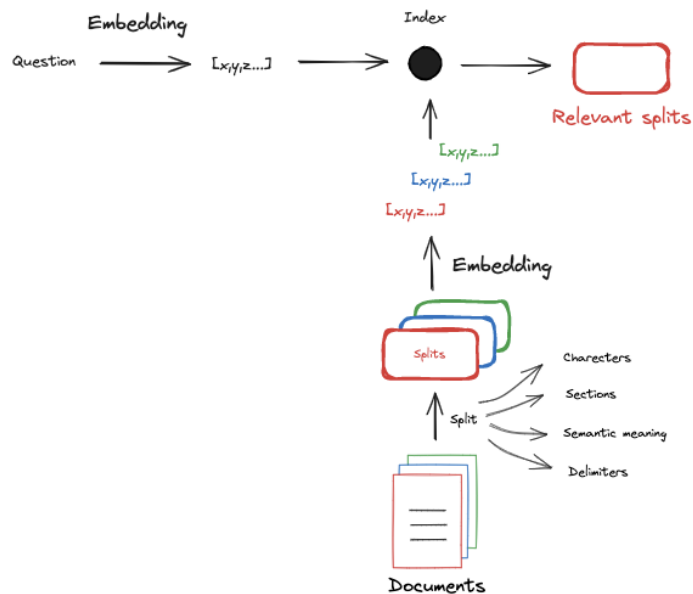


<https://simonwillison.net/2023/Oct/23/embeddings/>

<https://www.pinecone.io/learn/series/nlp/dense-vector-embeddings-nlp/>

https://cameronrwolfe.substack.com/p/the-basics-of-ai-powered-vector-search?utm_source=profile&utm_medium=reader2

Loading, splitting, and embedding

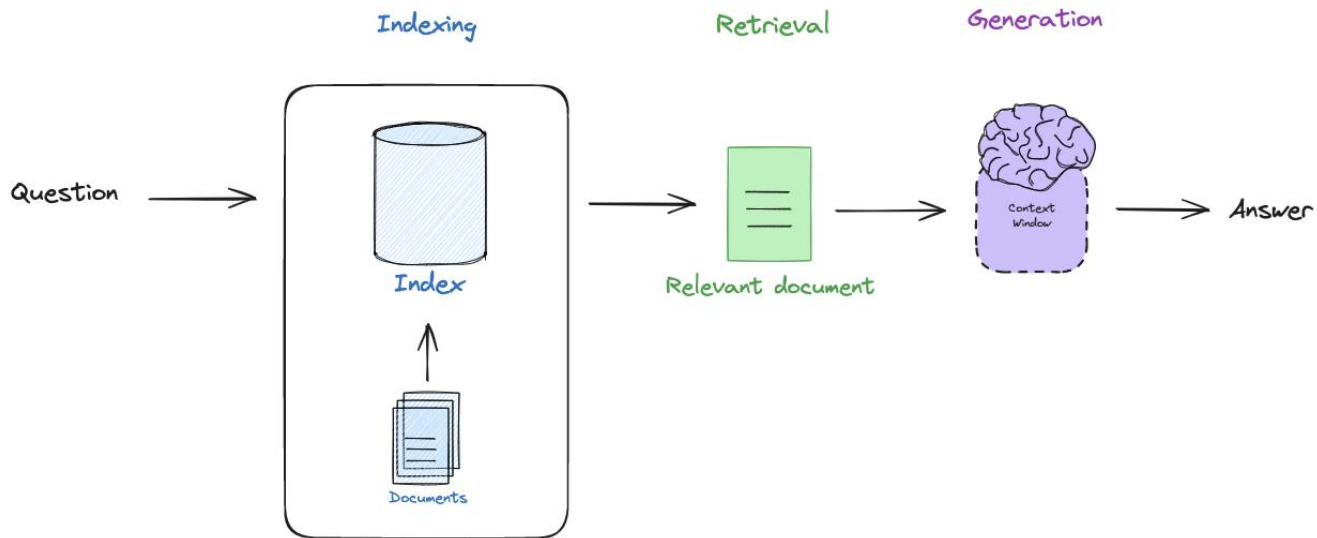


Code walk-through

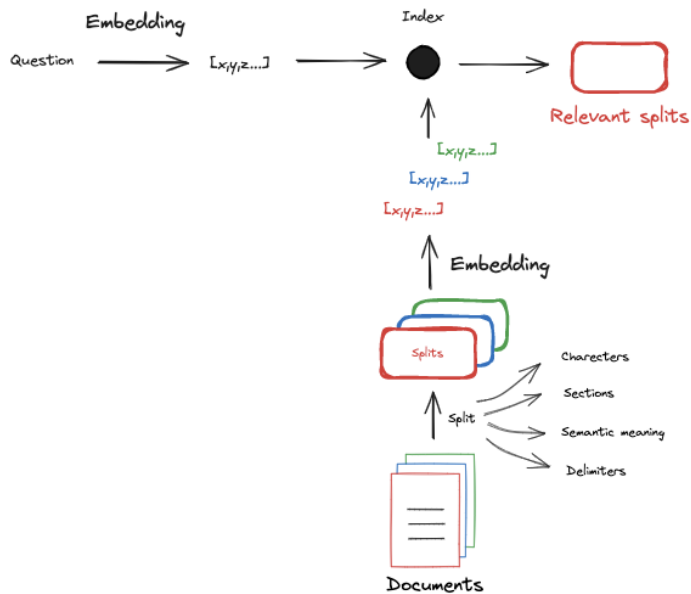
RAG from scratch: Retrieval

Lance Martin
Software Engineer, LangChain
[@RLanceMartin](https://twitter.com/RLanceMartin)

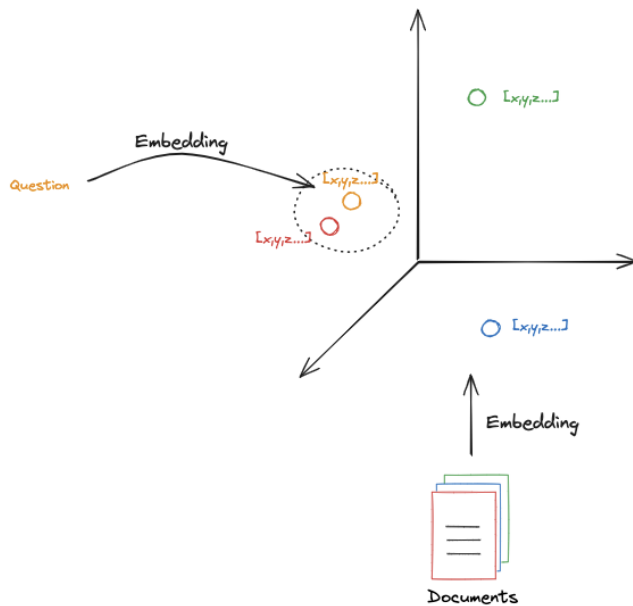
RAG motivation



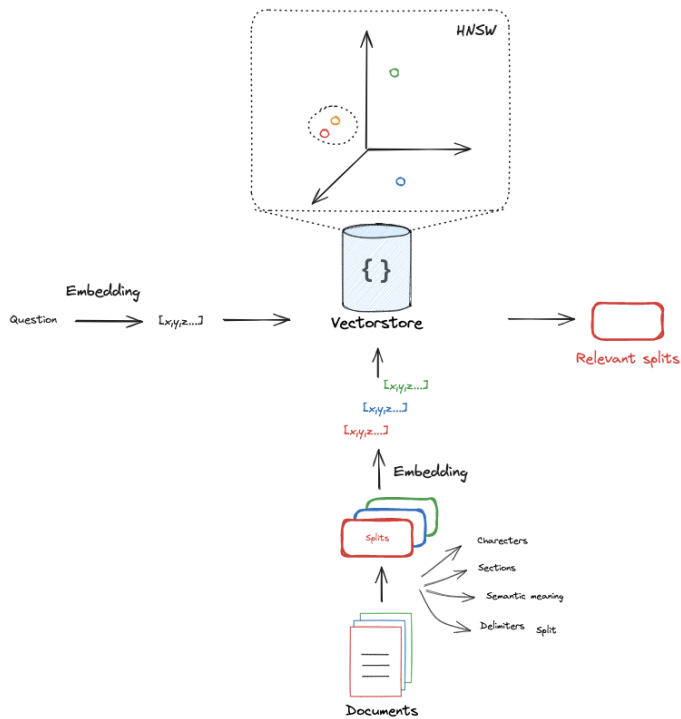
Index makes documents easy to retrieve



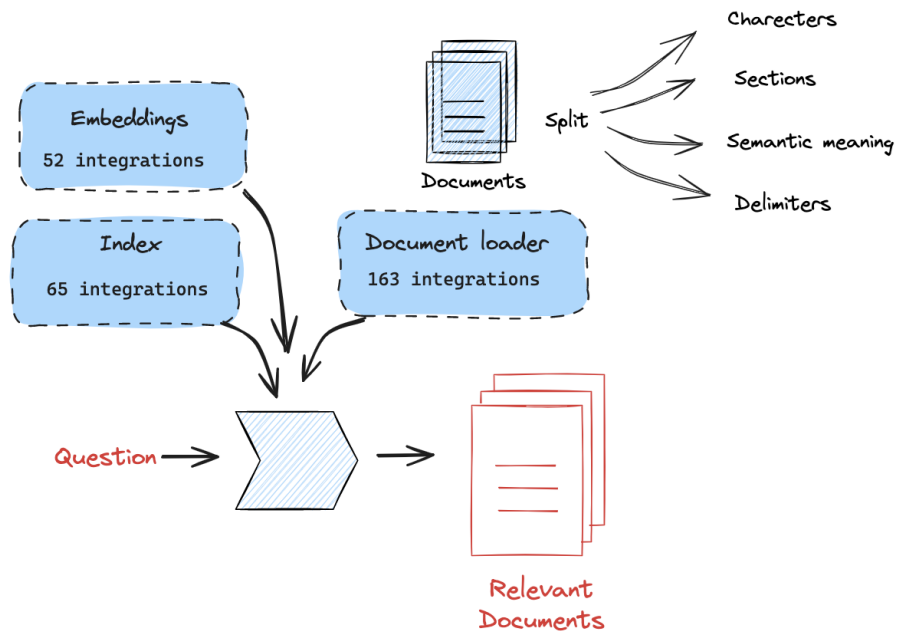
Retrieval powered via similarity search



Vectorstores implement this for you



LangChain has many integrations to support this

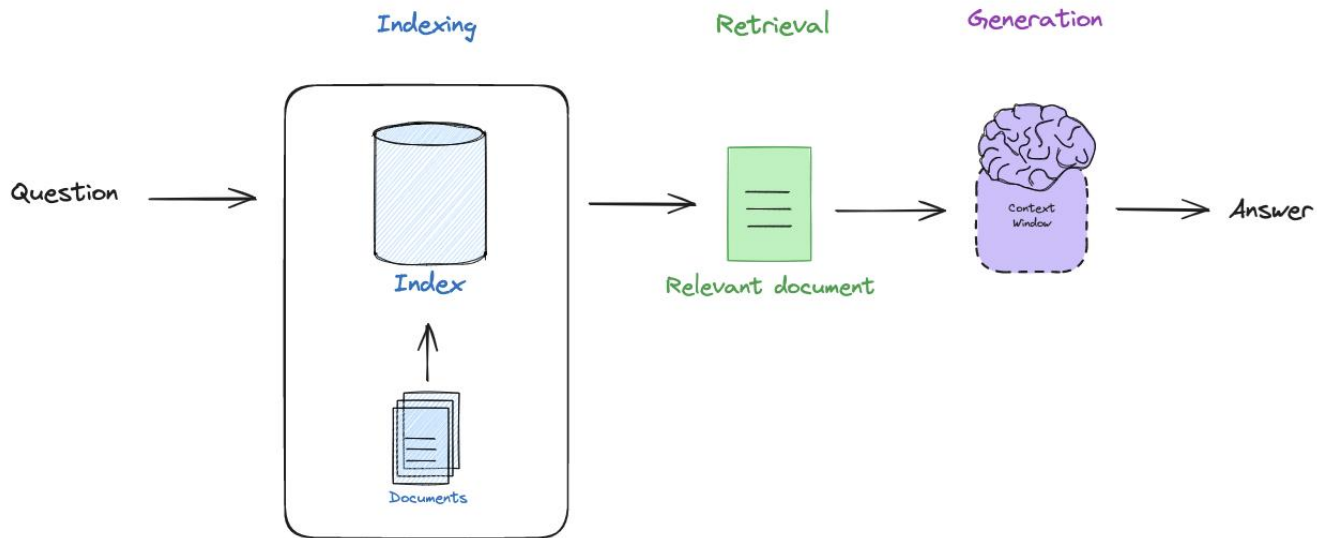


Code walk-through

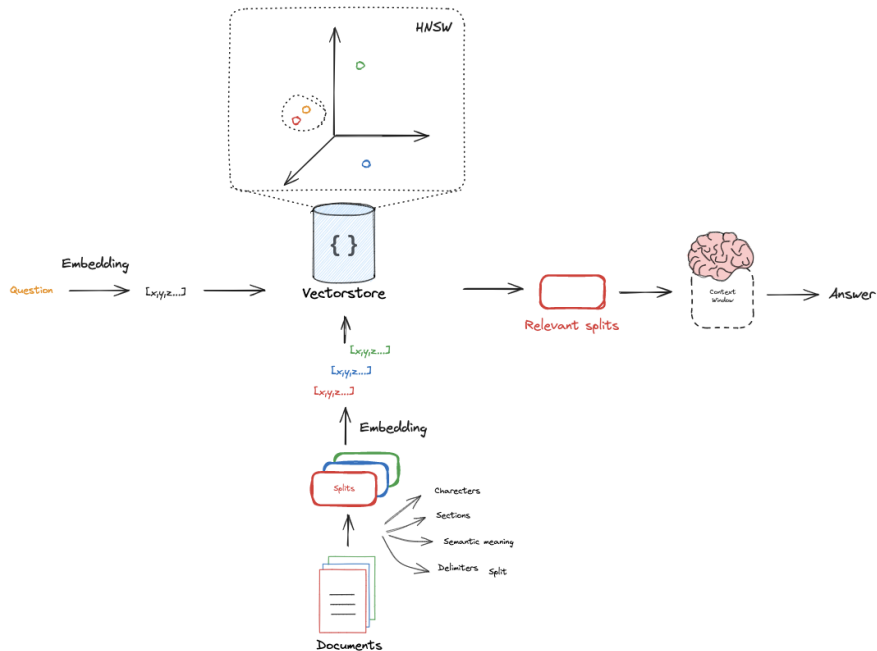
RAG from scratch: Generation

Lance Martin
Software Engineer, LangChain
[@RLanceMartin](https://twitter.com/RLanceMartin)

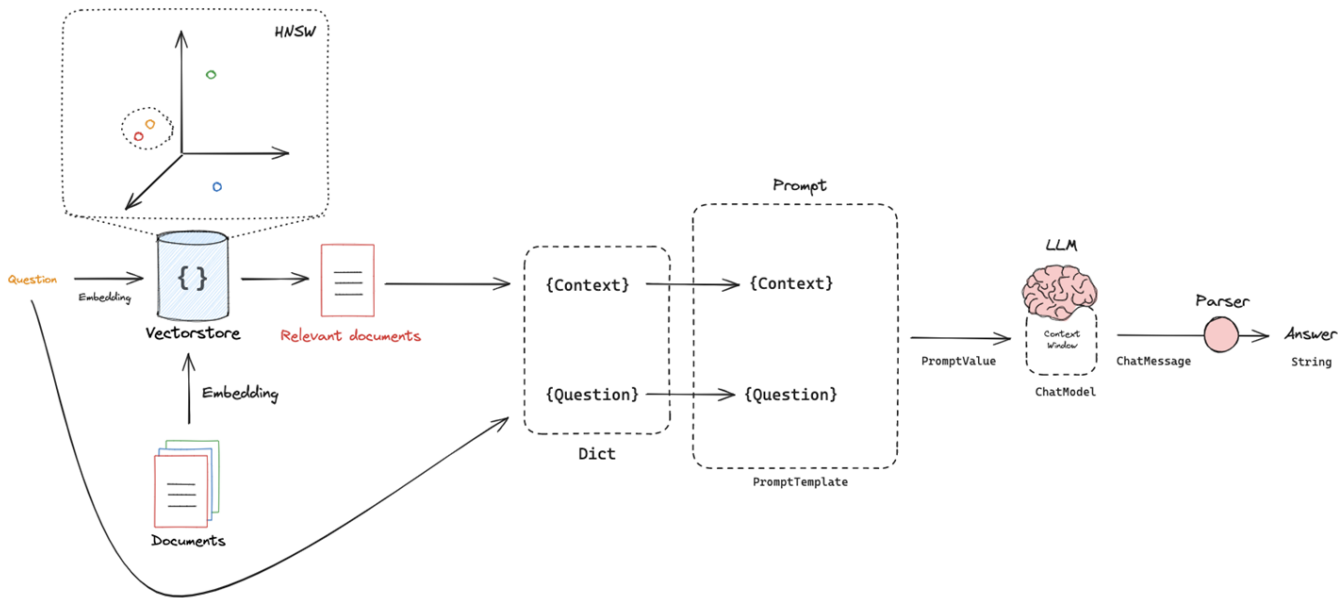
RAG motivation



Adding docs to context window



Connecting retrieval with LLMs via prompt



Code walk-through