

Retrieval Augmented Generation

Duarte Santos – 124376

Table of Contents

What is RAG?

RAG structure

Advantages and Limitations of RAG

RAG Use Cases

Conclusions

What is Retrieval Augmented Generation?

RAG is a framework, proposed by Patrick Lewis in 2020, that combines retrieval-based systems with generative AI models to improve the relevance and accuracy of generated content.

This approach leverages both explicit knowledge from retrieved documents and the implicit knowledge encoded in language models, resulting in more accurate and context-sensitive responses.

RAG Components



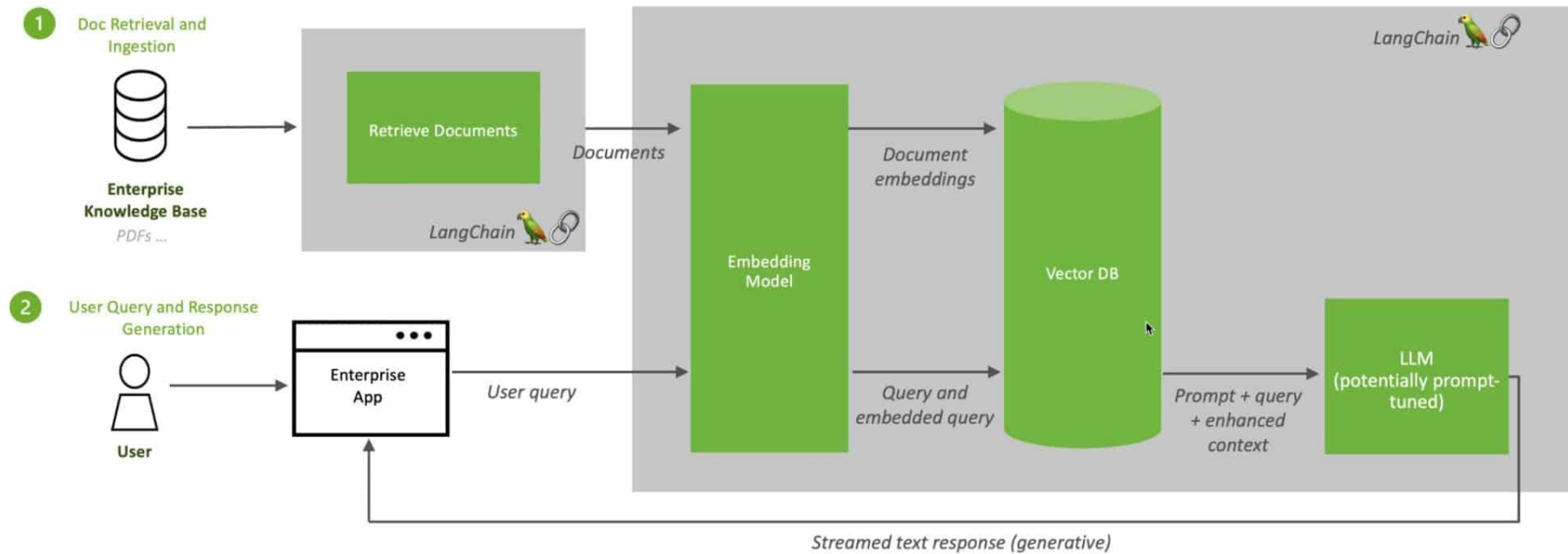
Retrieval: The system queries a knowledge base or external data source (e.g., databases, documents, or APIs) to fetch relevant context or information based on the input prompt



Generation: A generative model, such as GPT, uses the retrieved information to produce a coherent and contextually accurate response or output

How RAG works?

Retrieval Augmented Generation (RAG) Sequence Diagram



<https://blogs.nvidia.com/blog/what-is-retrieval-augmented-generation/>

RAG Example

Who were the candidates for the 2024 USA election and who won in November?



I might not have the latest updates. For the most complete and up-to-date information about the U.S. Election, please visit news sources like [The Associated Press](#) and [Reuters](#) or check with your state or local election authority.

What is the last year this GPT model was trained on?



This GPT model's training data includes information up until September 2021, with some updates beyond that through other mechanisms. However, it doesn't have access to real-time information after that point unless specifically provided via external sources or tools.



100

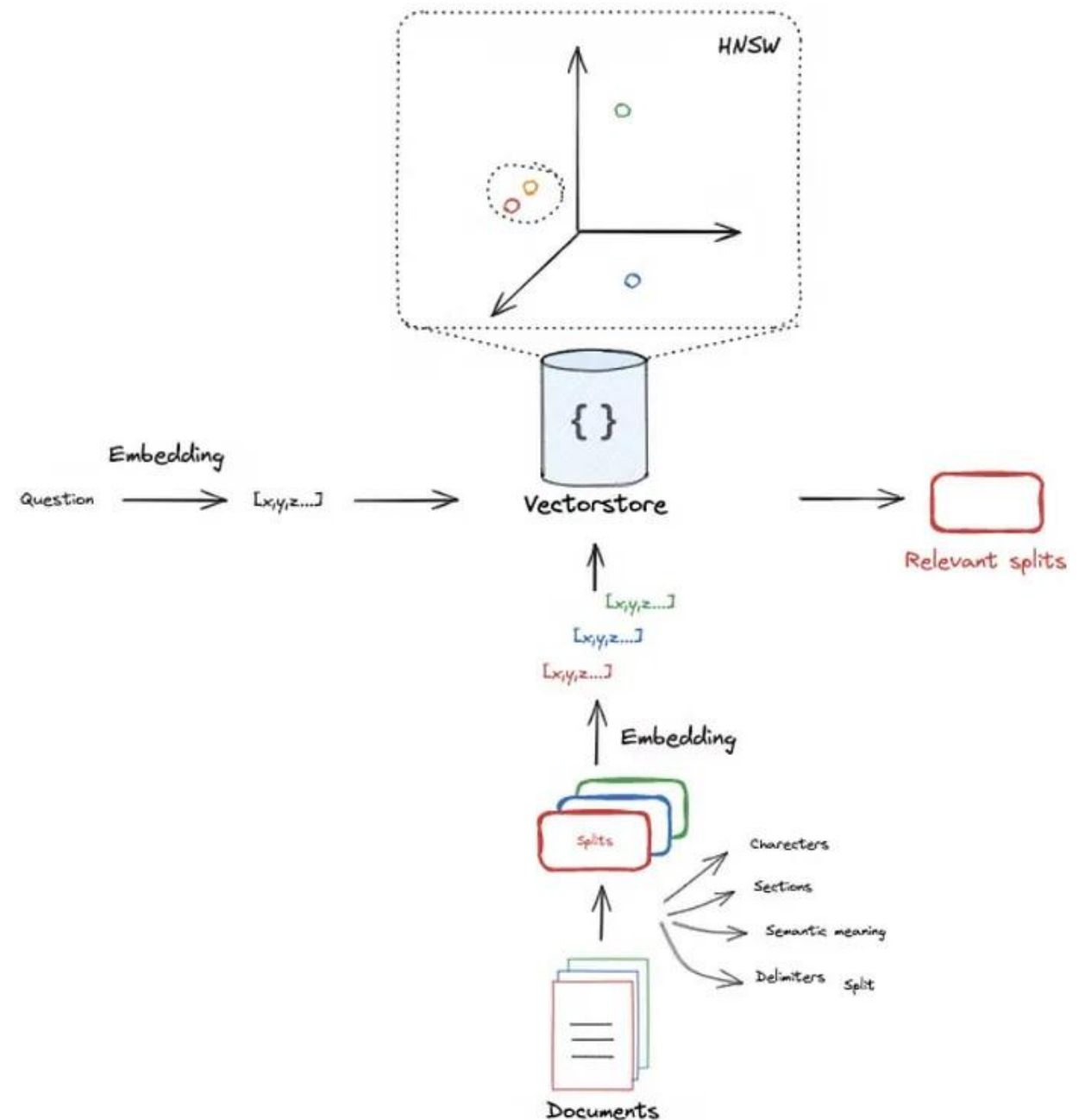
🔊 📄 👍 👎 ↺ ✓

Question: Who were the candidates for the 2024 USA election and who won in November?

08/01/2025

Embeddings improvements

Split the documents into smaller chunks to enable faster searching and more specific retrieval.



Advantages and Limitations of RAG

- Improved Accuracy
- Reduced LLMs Hallucinations
- Modularity
- Dynamic Knowledge Update
- Domain Adaptability

- Complex Architecture
- Latency
- Potential for Irrelevant Context
- Scalability of Retrieval

RAG Use Cases



CUSTOMER SUPPORT
CHATBOTS



SCIENTIFIC
RESEARCH



LEGAL ASSISTANCE



ENTERPRISE
DOCUMENT SEARCH

Future of RAG



Integration with Real-Time and Dynamic Knowledge



Multilingual retrieving and synthesizing evidence from diverse multilingual sources



Multimodal data (images, video, sound etc...) support



Smaller Models with External Knowledge

Conclusions

Retrieval Augmented Generation offers flexibility, adaptability, and efficiency, making it ideal for a wide range of use cases.

As research and technology progress, RAG systems will become more sophisticated, real-time, and domain-specific, opening new frontiers for AI applications



Questions?

