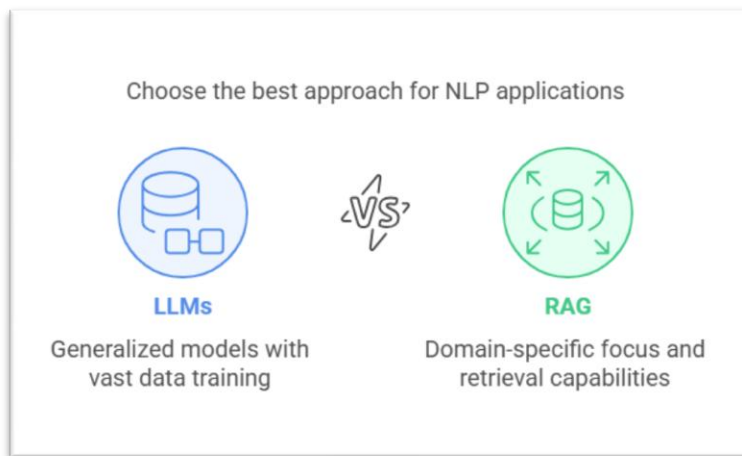# Unlocking the Power of RAG: Simple, Advanced, and Cache RAG
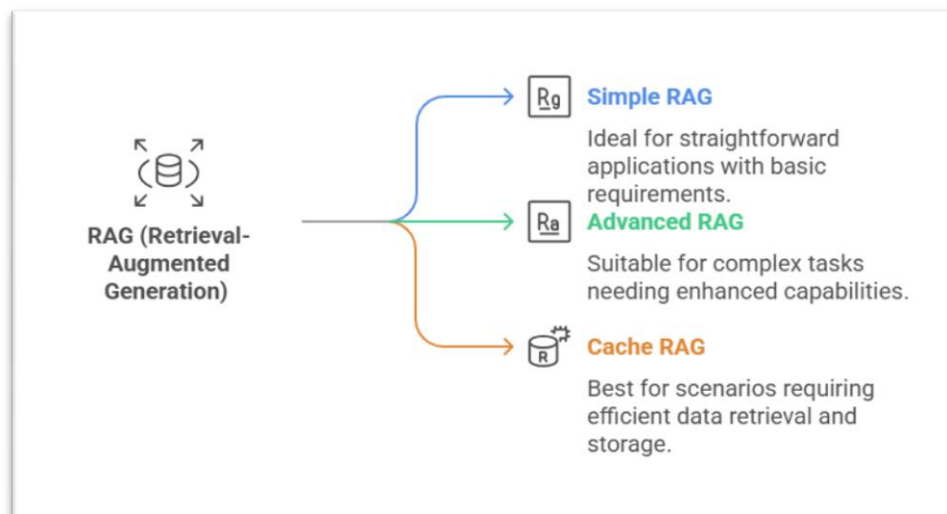
Everyone knows a buzzword: LLMs. But what are these LLMs? They form the foundation of RAG.

RAG stands for Retrieval-Augmented Generation. LLMs are models trained on vast amounts of data, making them generalized models. However, RAG is primarily designed to address domain-specific challenges.Although LLMs play a crucial role in developing NLP applications, they have certain limitations, such as hallucinations, a lack of domain-specific knowledge, no direct source attribution, and potentially outdated training data. In such scenarios, RAG plays a significant role in enhancing NLP applications.



## In This blog Let us go through 3 types of Rag:

1. Simple RAG
2. Advanced RAG
3. Cache RAG (CAG)
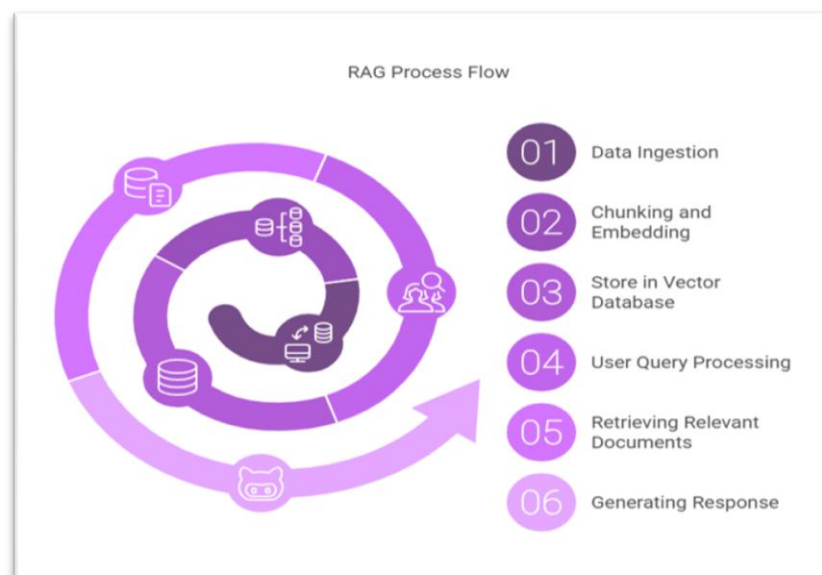4. HuggingFace models using Inference API

# 1. Simple Rag:

What is simple RAG? 🥲 Yes, it is the most basic version of RAG. In this structure, RAG operates with minimal capacity—receiving a user query, processing it using data stored in a vector database, and returning the required answerKey components include data ingestion, chunking, embedding generation, and vector storage, while retrieval strategies, prompt engineering, and LLM selection are crucial for performance. Effective RAG requires managing context window limitations and optimizing retrieval and generation processes. Ultimately, Simple RAG serves as a foundational step toward more advanced RAG architectures.

Before diving into the details of RAG, we need to understand some key concepts involved in building a RAG system:

1. Data Ingestion
2. Chunking and Embedding
3. Storing in a Vector Database
4. User Query Processing
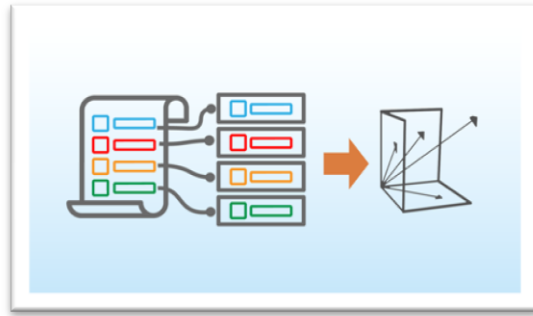5. Retrieving Relevant Documents
6. Generating a Response



RAG Process Flow

01 Data Ingestion
02 Chunking and Embedding
03 Store in Vector Database
04 User Query Processing
05 Retrieving Relevant Documents
06 Generating Response

# 1. Data Ingession:

This is the first and most crucial step in creating any RAG application because RAG primarily relies on data. Data is the key to any LLM model.

In this step, we gather the data that will be fed into the LLM. The data can come in various formats, such as text files, PDFs, Word documents, or website content. This data is then made accessible to the model for processing.
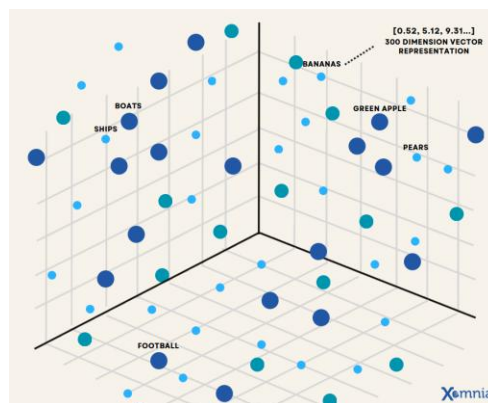
## 2. Chunking and embedding:



By looking at the above image, we can understand the concepts of chunking and embeddings. Chunking is the process of dividing data into smaller, manageable parts. For example, if the chunk size is set to 100, each chunk will contain 100 words or sentences, as specified. Another important parameter is overlap, which ensures semantic continuity. If the overlap value is 10, each new chunk will include 10 words from the previous chunk to preserve meaning. Once chunking is complete, these chunks are converted into vector embeddings, meaning they are transformed into numerical representations using techniques like cosine similarity. This ensures that semantically similar data points remain close to each other in the vector space.

## 3. Vector databases:

This step involves storing the embedding chunks in vector databases, which organize vector embeddings in a structured manner for future use. Vector databases facilitate efficient similarity searches, enabling quick retrieval of relevant information based on user queries. Different vector databases are suited for specific use cases, with some optimized for speed and scalability, while others prioritize seamless integration with AI frameworks. Prominent vector databases include Chroma DB, FAISS, Pinecone, and Redis.
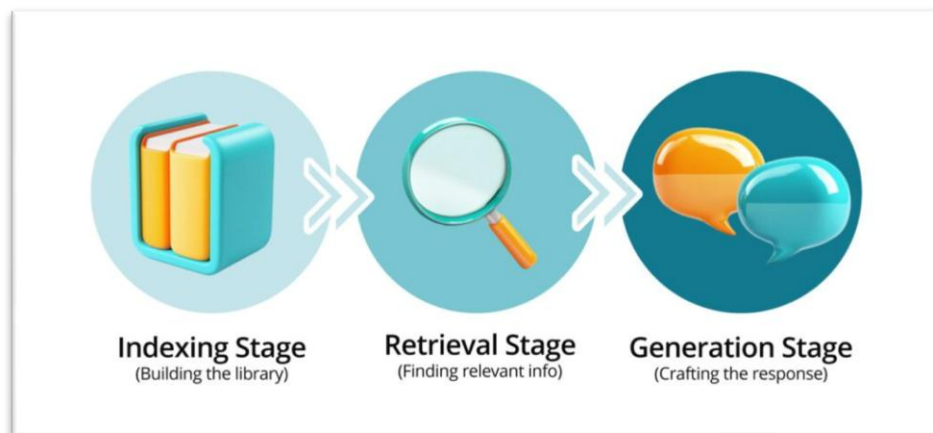


This is a small example of how the data is store in the database.

# 4. User query processing:

The next important stage is receiving and processing the user query. So far, we have covered data collection, chunking, embedding, and storing it in a vector database. Once these steps are completed, we obtain a query from the user related to the core concepts present in the stored data. The query is then embedded and converted into a numerical representation. After this, the crucial step of retrieving relevant documents begins.

# 5. Retrieving relevant documents:



**Indexing Stage**
(Building the library)

**Retrieval Stage**
(Finding relevant info)

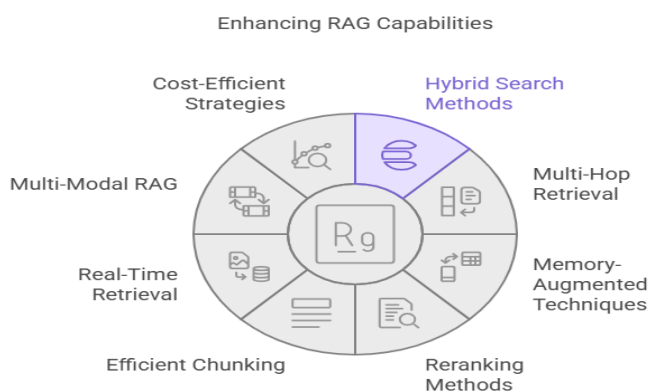**Generation Stage**
(Crafting the response)

This is the stage where we match the user query with the data stored in the database. Once the query is embedded, we compute its similarity with the stored vectors using cosine similarity. Two key parameters to consider here are **K** and **match probability**. **K** represents the top **K** number of documents that are most relevant to the user query. **Match probability** is a value ranging from 0 to 1, determining how strictly we filter the results. If the value is close to **1**, we retrieve only the vectors that are highly similar to the query, ensuring precise matching. As the value moves closer to **0**, we allow more flexibility, accepting vectors that may be slightly less similar but still relevant. Choosing the right values for these parameters plays a crucial role in optimizing retrieval accuracy and relevance.

# 6. Generating Response:

The final stage is generating the response using an LLM model, which can be from OpenAI, Gemini, Groq, or Cohere. These models are provided with the retrieved data from the vector database along with a prompt. By analyzing this data and the prompt, the LLM generates a response that is accurate and relevant. The quality of the response depends on the effectiveness of data retrieval and prompt engineering, ensuring that the output is clear, precise, and contextually appropriate.
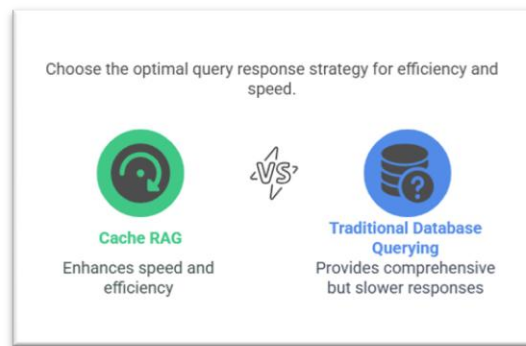
# 2. Advanced RAG:

Advanced Retrieval-Augmented Generation (RAG) enhances the traditional RAG pipeline by incorporating hybrid search methods, multi-hop retrieval, and memory-augmented techniques. Instead of relying solely on dense vector search, hybrid retrieval combines keyword-based (BM25) and semantic search for improved accuracy. Multi-hop retrieval allows the system to recursively query related documents, ensuring better context comprehension, while memory-augmented RAG retains past interactions for personalized and adaptive responses. Additionally, reranking methods like cross-encoder models refine the retrieved documents before passing them to the LLM, optimizing relevance. Efficient chunking techniques, such as hierarchical and dynamic chunking, further enhance retrieval by structuring text more effectively.



To scale and improve response efficiency, RAG integrates real-time retrieval, federated search, and streaming vector databases that support continuous updates. Multi-modal RAG extends the capability to handle images, audio, and video, using cross-modal embeddings for diverse content understanding. Cost-efficient strategies include approximate nearest neighbor (ANN) search, query caching, and on-device retrieval to minimize latency and API expenses. Furthermore, integrating external tools like Google Search APIs, code execution, and domain-specific fine-tuning refines the system for specialized applications. These advancements make RAG more robust, adaptable, and capable of handling complex, real-world information retrieval tasks with higher precision.
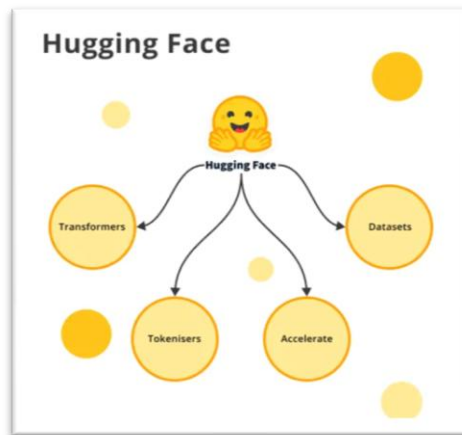
# 3. Cache RAG (CAG):

Cache RAG is an advanced concept in RAG. Instead of querying the database every time for an answer, it utilizes a temporary memory or cache. In this approach, previously asked user queries and their responses from the LLMs are stored. These cached responses are used when similar types of questions arise again. For example, if we have a database trained on data science concepts and a user asks the query, "What is data science?" the system will initially search the database for relevant information and provide an answer. However, in CAG, after generating the response, the system stores both the query and the response in the cache. If another user later asks a similar query, the system will first check the cache instead of searching the database, allowing for quicker retrieval.

Choose the optimal query response strategy for efficiency and speed.

Cache RAG
Enhances speed and efficiency

'VS'

Traditional Database Querying
Provides comprehensive but slower responses

This approach significantly improves response speed, as the system does not need to repeatedly query the database for frequently asked questions. By reducing the number of database lookups, it also enhances system efficiency and reduces computational costs. Additionally, caching helps in optimizing memory usage, ensuring that repeated queries do not add unnecessary processing load. However, one challenge with CAG is cache management—ensuring that outdated or irrelevant responses are refreshed periodically. Implementing a smart caching strategy, such as time-based expiration or update triggers, can further enhance the effectiveness of Cache RAG in real-world applications.

# 4. Hugging Face Models using Inference API:



Hugging Face

Transformers

Datasets

Tokenisers

Accelerate

Now let us discuss the last topic ie hugging face models and inference api. The Hugging Face Inference API allows users to deploy and run machine learning models without requiring local installations or high-end hardware. It provides access to a wide range of pre-trained models, including transformers for NLP, vision, and audio tasks. By making API calls to hosted models, users can perform tasks like text generation, sentiment analysis, and image classification with minimal setup. This approach ensures scalability, reduces infrastructure costs, and simplifies integration into applications.