

Here's a brief overview of the RAG pipeline:

1. **Retrieval:** This stage involves searching and retrieving relevant information from a knowledge source (e.g., documents, databases, or the internet) based on a user's query. Various techniques, such as vector-based search or sparse indexing, can be used to efficiently retrieve the most relevant information.
2. **Generation:** Once relevant information is retrieved, an LLM is used to generate an answer by considering both the user's query and the retrieved information. The LLM processes the retrieved information and combines it with its internal knowledge to produce a coherent and accurate answer.
3. **Re-ranking:** In some cases, multiple candidate answers may be generated by the LLM. A re-ranking step can be applied to score and sort the candidate answers based on their relevance, accuracy, or other criteria.

## **What are the frameworks and libraries available to implement RAG?**

Yes, there are several frameworks and libraries available that can help in implementing a RAG (Retrieval-Augmented Generation) pipeline.

These tools provide components for retrieval, integration with Large Language Models (LLMs), and generation of responses. Here are a few popular frameworks:

1. **Hugging Face Transformers:** Hugging Face offers a wide range of transformer-based models, including LLMs and retrieval models like Dense Passage Retrieval (DPR). You can combine these models to create a RAG pipeline for question answering or content generation tasks.
2. **LangChain:** LangChain is a framework for composable Language Model Development, which simplifies the integration of LLMs with other components, like retrieval models. It provides a flexible and modular approach to building RAG pipelines and other complex language model applications.
3. **Haystack:** Haystack is an open-source framework for building search systems using the latest deep learning models, including transformer-based LLMs. It supports various retrieval models and document stores, making it a versatile choice for implementing the retrieval part of a RAG pipeline.
4. **Deepset Haystack + Transformers + ONNX Runtime:** This combination of tools can be used to create an end-to-end RAG pipeline. Deepset Haystack can handle document processing and storage, Transformers provide the LLM for

generation, and ONNX Runtime allows for efficient and scalable deployment.