

Retrieval-Augmented Generation (RAG) is a technique in NLP that enhances the performance of language models by combining a retrieval system with a generation model. It retrieves relevant documents from a knowledge base and feeds them into a generative model to provide accurate and grounded responses.

FAISS (Facebook AI Similarity Search) is an open-source library that enables efficient similarity search and clustering of dense vectors, which is commonly used to implement the retrieval part of a RAG system.

FLAN-T5 is a fine-tuned T5 model developed by Google that can be used for various text generation tasks and performs well with few-shot learning. It's often used as the generator in RAG pipelines.

This document will serve as an example knowledge base for demonstrating a RAG pipeline using LangChain and Google Colab.