

I'll break down the code **step by step**, explain **each concept**, provide **theoretical background**, and include **an architectural diagram**.

Overview of the Code

The provided code implements **efficient document reranking** and **retrieval-augmented generation (RAG)** using **FlashRank**, **FAISS**, and **LangChain**.

It focuses on:

1. **FlashRanker-based Reranking**: Reorders retrieved documents based on relevance.
2. **Embedding-based Retrieval with FAISS**: Uses vector search to fetch relevant documents.
3. **Contextual Compression with FlashRankRerank**: Compresses retrieved data before passing it to an LLM.
4. **Retrieval-Augmented Generation (RAG)**: Uses retrieved documents to enhance LLM responses.

Step-by-Step Code Explanation

Step 1: Install Required Libraries

```
python
```

```
CopyEdit
```

```
!pip install flashrank langchain_community langchain_openai faiss-cpu
```

Installs necessary libraries:

- `flashrank` → Lightweight reranking model.
- `langchain_community` → LangChain utilities.
- `langchain_openai` → OpenAI integration for embeddings & chat models.
- `faiss-cpu` → Efficient vector database for retrieval.

Step 2: Helper Function to Print Documents

```
python
CopyEdit
def pretty_print_docs(docs):
    print(
        f"\n{'-' * 100}\n".join(
            [
                f"Document {i+1}:\n\n{d.page_content}\nMetadata:
{d.metadata}"
                for i, d in enumerate(docs)
            ]
        )
    )
```

✅ **Formats retrieved documents for display.**

Step 3: Define Query and Passages

```
python
CopyEdit
query = "How to speedup LLMs?"
passages = [
    {"id":1, "text":"Introduce *lookahead decoding*...", "meta":
{"additional": "info1"}},
    {"id":2, "text":"LLM inference efficiency...", "meta":
{"additional": "info2"}},
    {"id":3, "text":"There are many ways to increase LLM inference
throughput...", "meta": {"additional": "info3"}},
    {"id":4, "text":"Ever want to make your LLM inference go
brrrrrr...", "meta": {"additional": "info4"}},
    {"id":5, "text":"vLLM is a fast and easy-to-use library...",
"meta": {"additional": "info5"}}
]
```

- ✅ Defines **sample text passages** that describe different LLM speed-up techniques.
- ✅ **Each passage has metadata** for additional context.

Step 4: Define the FlashRanker Function

```
python
CopyEdit
from flashrank.Ranker import Ranker, RerankRequest

def get_result(query, passages, choice):
    if choice == "Nano":
        ranker = Ranker()
    elif choice == "Small":
        ranker = Ranker(model_name="ms-marco-MiniLM-L-12-v2",
cache_dir="/opt")
    elif choice == "Medium":
        ranker = Ranker(model_name="rank-T5-flan", cache_dir="/opt")
    elif choice == "Large":
        ranker = Ranker(model_name="ms-marco-MultiBERT-L-12",
cache_dir="/opt")

    rerankrequest = RerankRequest(query=query, passages=passages)
    results = ranker.rerank(rerankrequest)
    print(results)

    return results
```

- ✅ Defines a function `get_result()` that uses **FlashRanker to reorder documents**.
- ✅ **Chooses between different ranking models** based on user selection:
 - **Nano** → Fastest (~4MB)
 - **Small** → More accurate (~34MB)
 - **Medium** → Best zero-shot ranking (~110MB)
 - **Large** → Supports **100+ languages** (~150MB) ✅ Uses `RerankRequest()` to rank documents based on relevance to the query.

Step 5: Execute FlashRanker with Different Models

```
python
CopyEdit
%%time
get_result(query, passages, "Nano")

%%time
get_result(query, passages, "Small")

%%time
get_result(query, passages, "Medium")
```

✅ **Runs reranking for different models** and measures execution time.

Step 6: Load Documents & Split for Processing

```
python
CopyEdit
from langchain_community.document_loaders import TextLoader
from langchain_text_splitters import RecursiveCharacterTextSplitter

documents = TextLoader("/content/state_of_the_union.txt").load()
text_splitter = RecursiveCharacterTextSplitter(chunk_size=500,
chunk_overlap=100)
texts = text_splitter.split_documents(documents)
```

✅ **Loads a large document (state_of_the_union.txt)** and splits it into smaller chunks (**500 characters, with 100 overlapping characters**) using **RecursiveCharacterTextSplitter**.

Step 7: Assign Metadata to Documents

```
python
```

CopyEdit

```
for id, text in enumerate(texts):  
    text.metadata["id"] = id
```

- ✅ **Assigns unique IDs** to each document chunk.

Step 8: Generate Embeddings & Store in FAISS

python

CopyEdit

```
from langchain_community.embeddings import OpenAIEmbeddings  
from langchain_community.vectorstores import FAISS
```

```
embedding = OpenAIEmbeddings(model="text-embedding-ada-002")  
retriever = FAISS.from_documents(texts,  
embedding).as_retriever(search_kwargs={"k": 10})
```

- ✅ **Creates embeddings** for document chunks using **OpenAI's text-embedding-ada-002** model.

- ✅ **Stores embeddings in FAISS**, enabling **fast similarity search**.

- ✅ Retrieves **top-10** relevant documents for a query.

Step 9: Retrieve Documents Based on Query

python

CopyEdit

```
query = "What did the president say about Ketanji Brown Jackson"  
docs = retriever.invoke(query)  
pretty_print_docs(docs)
```

- ✅ **Searches FAISS** for relevant **State of the Union** speech snippets.
- ✅ Prints retrieved results.

Step 10: Compress Retrieved Documents with FlashRank

```
python
CopyEdit
from langchain.retrievers import ContextualCompressionRetriever
from langchain.retrievers.document_compressors import FlashrankRerank
from langchain_openai import ChatOpenAI

llm = ChatOpenAI(temperature=0)
compressor = FlashrankRerank()
compression_retriever =
ContextualCompressionRetriever(base_compressor=compressor,
base_retriever=retriever)

compressed_docs = compression_retriever.invoke("What did the president
say about Ketanji Jackson Brown")
pretty_print_docs(compressed_docs)
```

- ✅ Uses `FlashrankRerank()` **to further filter** retrieved documents.
- ✅ **Reduces hallucinations** by ensuring only highly relevant documents are passed to the LLM.

Step 11: Generate Final Answer with Retrieval-Augmented Generation (RAG)

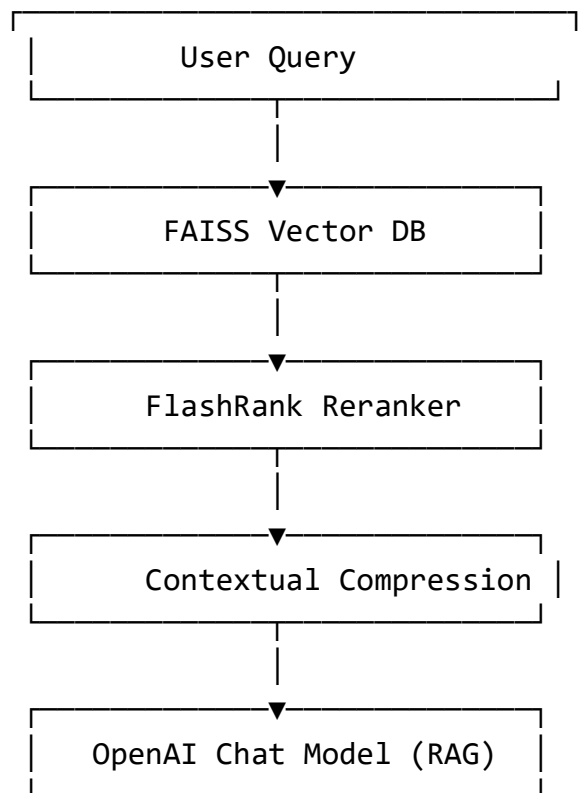
```
python
CopyEdit
from langchain.chains import RetrievalQA

chain = RetrievalQA.from_chain_type(llm=llm,
retriever=compression_retriever)
chain.invoke(query)
```

- ✓ Uses **retrieved and compressed documents** as context for **LLM-generated responses**.
- ✓ **Ensures factual accuracy** by limiting model output to **retrieved context**.

Architecture Diagram

graphql
CopyEdit



Missing Concepts (Improvements)

- ✓ **Use Hybrid Retrieval** (BM25 + Embeddings).
- ✓ **Experiment with Other Rerankers** (ColBERT, RankT5).
- ✓ **Implement Query Expansion** for better search results.

Final Thoughts

This code **efficiently retrieves, ranks, and filters** documents before feeding them to an LLM, ensuring **speed, accuracy, and factual correctness**. 🚀

```
#####DeepSeek#####
```

1. Line-by-Line Explanation

Installation & Setup

```
python
Copy
!pip install flashrank
```

- Installs the `flashrank` library for lightweight re-ranking of search results.

Helper Function

```
python
Copy
def pretty_print_docs(docs):
    print(
        f"\n{'-' * 100}\n".join(
            [
                f"Document {i+1}:\n\n{d.page_content}\nMetadata: {d.metadata}"
                for i, d in enumerate(docs)
            ]
        )
    )
```

- **Purpose:** Formats and prints documents with their metadata.
- **Mechanics:**
 - Iterates over a list of documents.
 - Separates documents with --- lines.
 - Includes page content and metadata for each document.

Query & Passages

```
python
Copy
query = "How to speedup LLMs?"
passages = [ ... ]
```

- **Query:** A question about speeding up large language models (LLMs).
- **Passages:** A list of 5 example documents (with id, text, and metadata) to be re-ranked.

FlashRank Initialization

```
python
Copy
from flashrank.Ranker import Ranker, RerankRequest

def get_result(query, passages, choice):
    if choice == "Nano":
```

```

    ranker = Ranker() # Uses default model (ms-marco-TinyBERT-L-
2-v2)
    elif choice == "Small":
        ranker = Ranker(model_name="ms-marco-MiniLM-L-12-v2",
cache_dir="/opt")
    elif choice == "Medium":
        ranker = Ranker(model_name="rank-T5-flan", cache_dir="/opt")
    elif choice == "Large":
        ranker = Ranker(model_name="ms-marco-MultiBERT-L-12",
cache_dir="/opt")

rerankrequest = RerankRequest(query=query, passages=passages)
results = ranker.rerank(rerankrequest)
print(results)
return results

```

- **Purpose:** Reranks passages using a selected FlashRank model.
- **Mechanics:**
 - Initializes a Ranker with a model based on the choice (Nano/Small/Medium/Large).
 - Creates a RerankRequest object with the query and passages.
 - Executes reranking and returns results.

Timing Execution

```

python
Copy
%%time
get_result(query, passages, "Nano")

```

- **Purpose:** Measures the execution time of reranking with the Nano model.

LangChain Integration

```

python
Copy
!pip install langchain_community langchain_openai
from google.colab import userdata

```

```
OPENAI_API_KEY = userdata.get('OPENAI_API_KEY')
os.environ["OPENAI_API_KEY"] = OPENAI_API_KEY
```

- Installs LangChain libraries and sets up OpenAI API key.

Document Loading & Splitting

```
python
Copy
from langchain_community.document_loaders import TextLoader
from langchain_text_splitters import RecursiveCharacterTextSplitter

documents = TextLoader("/content/state_of_the_union.txt").load()
text_splitter = RecursiveCharacterTextSplitter(chunk_size=500,
chunk_overlap=100)
texts = text_splitter.split_documents(documents)
for id, text in enumerate(texts):
    text.metadata["id"] = id
```

- **Purpose:** Loads a text file and splits it into chunks for processing.
- **Mechanics:**
 - Uses TextLoader to load a document.
 - Splits text into 500-token chunks with 100-token overlap.
 - Adds an id to each chunk's metadata.

Vector Store & Retriever

```
python
Copy
from langchain_community.embeddings import OpenAIEmbeddings
from langchain_community.vectorstores import FAISS

embedding = OpenAIEmbeddings(model="text-embedding-ada-002")
retriever = FAISS.from_documents(texts,
embedding).as_retriever(search_kwargs={"k": 10})
```

- **Purpose:** Creates a vector store for semantic search.
- **Mechanics:**

- Uses OpenAI embeddings (text-embedding-ada-002) to convert text to vectors.
- Builds a FAISS index for efficient similarity search.
- Retriever fetches the top 10 documents for a query.

Re-Ranking with FlashRank

```
python
Copy
from langchain.retrievers import ContextualCompressionRetriever
from langchain.retrievers.document_compressors import FlashrankRerank

compressor = FlashrankRerank()
compression_retriever = ContextualCompressionRetriever(
    base_compressor=compressor, base_retriever=retriever
)
compressed_docs = compression_retriever.invoke("What did the president
say about Ketanji Jackson Brown")
```

- **Purpose:** Combines retrieval and re-ranking to improve result quality.
- **Mechanics:**
 - FlashrankRerank reorders documents based on relevance.
 - ContextualCompressionRetriever wraps the retriever and compressor.

Question Answering

```
python
Copy
from langchain.chains import RetrievalQA

chain = RetrievalQA.from_chain_type(llm=llm,
retriever=compression_retriever)
chain.invoke(query)
```

- **Purpose:** Answers questions using the compressed/reranked documents.
- **Mechanics:**
 - Uses ChatOpenAI (LLM) to generate answers.
 - Combines retrieval, re-ranking, and LLM inference.

2. Key Concepts

Re-Ranking

- **What:** Reordering search results to prioritize relevance.
- **Why:** Initial retrievers (e.g., FAISS) use vector similarity, which may miss contextual relevance.
- **FlashRank:** Uses lightweight cross-encoder models to score query-passage pairs.

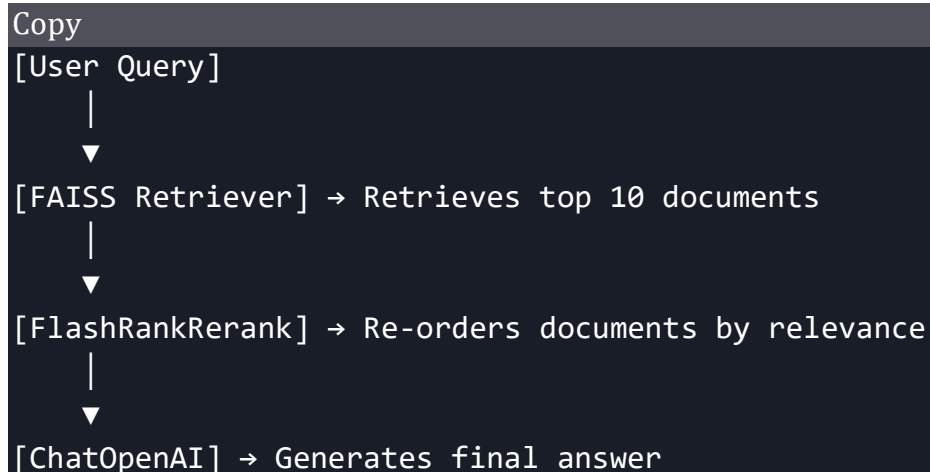
Cross-Encoders

- **Mechanism:** Process the query and document together (unlike bi-encoders).
- **Advantage:** Higher accuracy but slower than bi-encoders.
- **Models:**
 - **Nano:** TinyBERT (4MB, fastest).
 - **Small:** MiniLM (34MB, best accuracy).
 - **Medium:** T5-FLAN (110MB, zero-shot).
 - **Large:** MultiBERT (150MB, multilingual).

LangChain Components

- **Retriever:** Fetches documents (FAISS here).
- **Compressor:** Improves results (FlashrankRerank here).
- **Chain:** Combines retrieval, compression, and LLM (RetrievalQA).

3. Architecture Diagram



4. Missing Points

1. **Model Details:**

- a. FlashRank models are based on **cross-encoders**, which are slower but more accurate than bi-encoders.
- b. rank-T5-flan (Medium) is a text-to-text model, while others are BERT-based.

2. **Use Case:**

- a. Designed for serverless deployments (AWS Lambda, Vercel) due to low memory usage.

3. **Performance Tradeoffs:**

- a. Nano is 4MB but less accurate; Small is 34MB with best accuracy.

4. **Alternatives:**

- a. Compare with other rerankers (Cohere, SentenceTransformers).

5. **Quantitative Results:**

- a. Missing benchmarks (e.g., latency vs. accuracy for each model).

6. **Error Handling:**

- a. No code for handling API failures or model loading errors.

7. **Cache Directory:**

- a. `cache_dir="/opt"` is specific to deployment environments (e.g., Docker/cloud VMs).