**Introduction**

I will be using a Retrieval-Augmented Generation (RAG) model, a hybrid framework that combines information retrieval with natural language generation to extract relevant information from PDFs. The RAG model works by first retrieving the most relevant documents or pieces of text from a database (in this case, the PDFs) using a retriever model, typically based on dense or sparse vector search methods. Once the information is retrieved, a generator model, usually a pre-trained transformer or LLM, generates a coherent response based on the retrieved content. This approach allows the model to access external knowledge from the PDFs and generate accurate, contextually rich responses to complex queries.

The retriever is responsible for finding relevant information from a large collection of documents or data. In the case of processing PDFs, the retriever searches through the text to identify and retrieve passages that are most related to the query. To achieve this, the text is first transformed into high-dimensional vectors through a process called embedding, where each piece of text is mapped to a numerical representation that captures its meaning and context. This is done using embedding models like BERT or Sentence Transformers.

Once these vectors are created, they are stored in a vector store. A vector store is a specialized database designed to efficiently store, manage, and search these high-dimensional vectors. When a query is input, the retriever converts the query into a vector and searches the vector store to find the most similar vectors (pieces of text) based on mathematical distance measures like cosine similarity. The closest matches are retrieved as the most relevant documents or passages.

After retrieving the relevant information, the Large Language Model (LLM) comes into play. This LLM, typically a transformer-based model, processes the retrieved passages and combines them with the original query to generate a final response. The LLM is designed to understand context, integrate information from multiple sources, and generate coherent and meaningful text based on the retrieved content. This step ensures that the generated response is both accurate and contextually aligned with the query, providing a well-rounded answer that draws on the external knowledge contained in the PDFs.

**Methodology**

Using the Langchain framework for the RAG process provides several key advantages that streamline and enhance the workflow for retrieval-augmented generation. Here's why I chose Langchain for the  RAG implementation:

1. Seamless Integration of Components: Langchain is designed to make it easy to combine various components needed in a RAG pipeline, such as retrievers, vector stores, and large language models (LLMs). Langchain allows you to focus on the higher-level logic rather than worrying about integrating each piece manually.
2. Modularity and Flexibility: Langchain's modular structure allows you to customize each step of the RAG pipeline. You can choose from various text splitters, retrievers, embeddings, etc. Langchain supports a wide range of vector stores (like Pinecone, FAISS, Chroma) and embedding models (OpenAI, Hugging Face models), making it easier to work with whichever vector storage solution suits your application best.
3. Chaining Capabilities: Langchain excels at chaining multiple steps together, which is crucial for the RAG process. For example, the framework allows you to create chains that involve splitting documents into

chunks, embedding them, storing them in a vector database, retrieving relevant chunks, and then generating answers using an LLM.

4. Pre-built Tools and Utilities: Langchain comes with several pre-built utilities, such as text splitters, retrievers, and document loaders.
5. Document Loading and Parsing: Handling different types of documents (PDFs, Word documents, HTML pages, etc.) is crucial for RAG-based tasks. Langchain provides easy-to-use document loaders that simplify the process of ingesting documents from multiple formats.
6. Community and Ecosystem: Langchain has a growing community and is continually expanding its features with support for new models and tools.

I chose Hugging Face because it offers a wide variety of pre-trained models for both embeddings and LLMs, allowing easy access to state-of-the-art solutions. Its API simplifies integration into your RAG pipeline, reducing the complexity of managing models. Additionally, Hugging Face ensures scalability and provides strong community support for ongoing improvements and optimizations.

Some important tools used and their description:

# Split the documetns

```
def text_splits(data, size: int, overlap: int): text_splitter = RecursiveCharacterTextSplitter(chunk_size=size, chunk_overlap=overlap) chunks = text_splitter.split_documents(data) return chunks
```

The function takes in data (your document collection), a size (the number of characters for each chunk), and an overlap (how much the chunks should overlap). It uses the RecursiveCharacterTextSplitter to split large documents into smaller, manageable chunks of text that have overlapping regions. It then returns these chunks for further processing.

# Load the LLM

```
llm = HuggingFaceHub( repo_id="mistralai/Mistral-7B-Instruct-v0.3", # Specify the LLM model_kwargs=dict(max_new_tokens=1024, temperature=0.1, verbose=False), # LLM parameters huggingfacehub_api_token="api_key" # HuggingFace API token )
```

The HuggingFaceHub is used to load the Mistral-7B-Instruct-v0.3 model from Hugging Face's model repository.The model_kwargs parameter specifies how the model will behave during inference: max_new_tokens =1024: Limits the number of tokens generated in a single response. temperature=0.1: Controls the randomness of the output (lower values make the output more deterministic). verbose=False: Ensures minimal logging.

# QA chain setup for this document

```
qa_stuff = RetrievalQA.from_chain_type( llm=llm, chain_type="stuff", # Using the "stuff" chain type for simple Q&A retriever=retriever, verbose=False, return_source_documents=False, )
```

This function sets up and connects everything together in a Q&A chain using the RetrievalQA class. It specifies the LLM and a retriever to retrieve relevant document chunks, and it uses the "stuff" chain type, which simply combines all retrieved information for generating a response. The verbose flag is set to False to suppress logs,

and return_source_documents is also False, meaning only the answer is returned, not the documents themselves.

**Challenges**

The first challenge was attempting to run the model locally, which resulted in numerous failures and runtime errors due to the absence of a GPU, leading me to switch to Google Colab for its free GPU access. In Colab, I encountered difficulties finding and testing suitable embedding and LLM models that were effective yet lightweight enough to perform efficiently within the platform's constraints. Additionally, categorizing the PDFs into predefined categories posed challenges, as the process required careful consideration of the parameters. Lastly, for security reasons, I tried decrypting my API key but it was problematic since Colab necessitates hard-coding it, limiting my options for safe usage and increasing the risk of exposure.

**Conclusion**

Despite these efforts to implement a RAG model, the results have not been perfect. I faced challenges such as incorrect categorization and some irrelevant information in the title and authors, which highlight areas for improvement. Also, it is noted in the directions, that I should be able to run this offline, something which is not true. Because, I need to download my models from HuggingFace and also use a personal API key for this. Once downloaded, the script is considered that it runs offline. Finally, there is potential for the RAG model, as I believe that with further fine-tuning and testing, the model's performance and accuracy can be enhanced. With only a couple of days work, I was able to bring these positive results.