# Project Report:

# HelpMate AI Project

**Kiran S. Dalmiya**
August 28, 2024

# 1. Project Objectives

The aim of the Mr. Helpmate AI project is to undertake the development of a generative search system that is able to answer questions based on the document of a life insurance policy. It uses Retrieval-Augmented Generation (RAG) technology, which blends searching the policy document with AI-generated answers. This system first searches for relevant information in the document and then uses AI to generate precise and contextually relevant responses based on that information.

# 2. Design of the Project

## 2.1. Embedding Layer
**Document Processing & Cleaning:** The PDF of the life insurance policy was processed, extraneous metadata and sections not relevant to the topic like blank pages were removed. Footer was extracted and its data is stored in section columns.

**Strategy:** Text is divided into meaningful chunks in such a manner so that sentence structure is not destroyed during the chunking process, in turn keeping the contextual information in each chunk intact.

**Embedding:** Pre-processed chunks are embedded using a pre-trained model from the SentenceTransformers library fine-tuned for semantic understanding.

## 2.2. Search Layer
**Query Design:** The following were the three specific queries designed to test the system to search and find relevant fragments from the Policy document.
1. What is the procedure for filing the claim?
2. Are there any penalties for early withdrawal?
3.  What is the grace period for premium payment?

**Search Mechanism:** The stored embeddings of the embedded queries were matched against the embeddings of the documents stored in ChromaDB. A mechanism to store recent searches was implemented through a cache mechanism in the application to speed the process for similar queries.

**Re-ranking:** The relevant chunks were re-ranked using a cross-encoding model to ensure the most appropriate information was prioritized.

## 2.3. Generation Layer
**Prompt Engineering:** The fully articulated prompt was engineered to be very prescriptive and to have a few-shot example. This led the LLM on how to properly generate responses given that the top search results have been re-ranked.
**Response Generation:** The model received a prompt, now passed to an LLM (GPT-3.5-turbo), and final answer chunks were selected from those generated by the search layer and used.

## 3. Challenges Faced

**Chunking:** The most significant challenge was to maintain the proper balance between chunk size and meaningful segmentation. Too extensive a chunk offers a chance to lose specific details; on the other hand, if the chunk is small, the context may be lost.

**Model Selection:** The project required a way to optionally use alternative embedding models and trade off between accuracy and computational efficiency.

**Prompt Design:** Building a complete and clear prompt will be necessary to engage the generative model in generating responses that are contextually appropriate.

## 4. Lessons Learned

**Effective Chunking:** The proper chunking of text really affects the quality of embedding and search results; therefore, chunk boundaries must hold meaning.
**Embedding Model Impact:** The nature of the embedding model would be the sole distinction factor in the effectiveness of the retrieval process and thus would directly impact the answers' correctness. Here I used the text-embedding-ada-002 model to create embeddings for the text chunks from the life insurance policy document; this model helped me make sure that the system could effectively retrieve and generate accurate answers based on the policy document.
**Prompt Engineering Importance:** A well-crafted prompt can dramatically improve the generative model's output, particularly when dealing with complex or nuanced queries.

# 5. Screenshots & Results

Here are some of the screenshots compiled during the testing phase of the project:
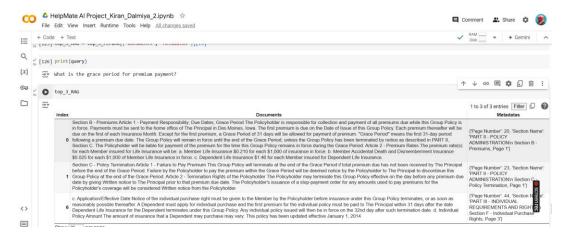
**Top 3 Search Layer Results:**

Screenshot 1: **Query 1 -** *What is the procedure for filing the claim?*



Screenshot 2: **Query 2 -** *Are there any penalties for early withdrawal?*

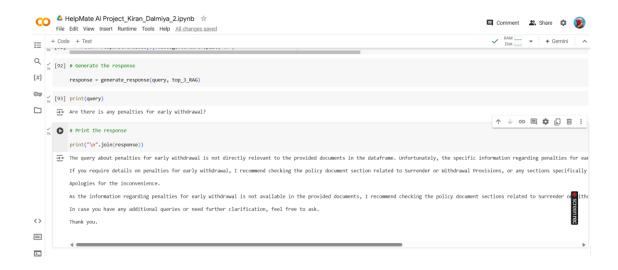Screenshot 3: **Query 3 - *What is the grace period for premium payment?***



**Final Generated Answer:**

*Screenshot 4:* **Query 1 - Final Answer**

*Screenshot 5:* **Query 2 - Final Answer**



*Screenshot 6:* **Query 3 - Final Answer**



These screenshots are indicating the system's capability of retrieving appropriate sections of documents and getting back responses with good structure.

## Conclusion

The Mr.Helpmate AI project developed a Retrieval-Augmented Generation system that integrates search and generative AI techniques. It also heavily focused on emphasizing text chunking, model selection, and prompt engineering to build a generative search system with efficacy and precision.

In the future, more advanced re-ranking techniques can be worked on, and multiple policy documents can be used to serve more generalized queries related to insurance.