# Report on the Working of the Movie Recommendation Chatbot Application

## OVERVIEW

The movie recommendation chatbot application combines a Large Language Model (LLM) with a vector database to create a personalized movie recommendation system using Retrieval-Augmented Generation (RAG). The primary components of this application include the main.py script, app.py script, and a dataset (movie_data.json). The application is designed to provide users with tailored movie recommendations based on their inputs and preferences.

## COMPONENTS.

**movie_data.json**:

This file contains the dataset of movies, each with attributes such as title, genre, and plot summary.

Example entry:

```
{
        "title": "Inside Out 2",
         "genre": "Animation",
         "plot": "Teenager Riley's mind headquarters is undergoing a sudden demolition to make room
         for something entirely unexpected: new Emotions! Joy, Sadness, Anger, Fear and Disgust,
         who've long been running a successful operation by all accounts, aren't sure how to feel when
         Anxiety shows up. And it looks like she's not alone."
}
```

**main.py**:

This script is responsible for handling the main logic of the application.

It includes the setup for the vector database and the integration with the LLM.

The script processes user inputs, retrieves relevant movie data from the vector database, and generates recommendations using the LLM.

Key functions:

load_data(): Loads movie data from movie_data.json.

create_vector_database(): Sets up the vector database with movie embeddings.

get_recommendations(): Retrieves movie recommendations based on user input.

**app.py**:

This script sets up the Flask web application for the chatbot interface.

It defines routes and handlers for user interactions.

Key routes:

/: Renders the homepage with the chatbot interface.

/recommend: Handles recommendation requests and returns responses generated by the LLM.

**movie_data.py:**

fetches the data from the TMDB database.

<div align="center">

**WORKFLOW**

</div>

**Data Loading and Vector Database Creation**:

The main.py script begins by loading the movie data from movie_data.json.

Each movie entry is processed to create vector embeddings, which are then stored in the vector database.

**User Interaction and Input Processing**:

The Flask application (app.py) serves the chatbot interface to users.

When a user provides input (e.g., requesting a movie recommendation), the input is sent to the server via an HTTP request.

**Retrieval-Augmented Generation (RAG)**:

**Retrieval**: The get_recommendations() function in main.py processes the user input by querying the vector database to find relevant movie entries. This involves generating vector embeddings for the user query and finding the nearest neighbors in the vector space, which represent the most relevant movies.

**Augmented Generation**: The retrieved movie data is then fed into the LLM, which uses this information to generate a coherent and contextually appropriate recommendation response. This combination of retrieval and generation ensures that the recommendations are both accurate and well-articulated.

**Response Delivery**:

The generated recommendations are sent back to the user through the Flask application.

The user receives the recommendations in the chatbot interface.

Example Interaction

1. **User Input**:

User: "movies where it is related to family"

2. **Processing**:

The input is processed to extract keywords and preferences.

The vector database is queried to find movies matching the criteria.

3. **Retrieval-Augmented Generation**:

   **Retrieval**: The query is transformed into vector embeddings and compared against the movie database to find relevant matches.

   **Augmented Generation**: The LLM generates a response using the retrieved movie data to craft a personalized recommendation.

   Example response:

   Avatar: The Way of Water (Science Fiction) - Set more than a decade after the events of the first film, learn the story of the Sully family (Jake, Neytiri, and their kids), the trouble that follows them, the lengths they go to keep each other safe, the battles they fight to stay alive, and the tragedies they endure.

   Sing and Don't Cry (Comedy) - A matriarchal family runs a taco food truck, navigating between artistic aspirations and daily frustrations while raising a member's daughter.

## CHALLENGES FACED

**Data Preprocessing**:

Ensuring the dataset is clean and structured correctly for vector embedding creation was time-consuming.

Handling missing or inconsistent data entries in the movie_data.json file required additional validation checks.

**Vector Database Integration**:

Setting up and optimizing the vector database to handle real-time queries efficiently was challenging.

Ensuring the database could scale with an increasing number of movie entries while maintaining fast retrieval times was essential.

**Embedding Generation**:

Generating accurate and meaningful vector embeddings for user queries and movie entries was crucial for effective recommendations.

Balancing the computational cost of embedding generation with the need for real-time performance required fine-tuning.

**LLM Integration**:

Integrating the LLM to work seamlessly with the retrieval system required careful handling of input and output data formats.

Ensuring the LLM generated contextually relevant and coherent responses based on the retrieved movie data was challenging.

**User Interaction**:

Designing a user-friendly chatbot interface that effectively captured user preferences and queries was critical.

Handling diverse and sometimes ambiguous user inputs required robust natural language processing capabilities.

## CONCLUSION

The movie recommendation chatbot application leverages advanced technologies to deliver personalized movie suggestions through Retrieval-Augmented Generation (RAG). By combining an LLM with a vector database, the application efficiently processes user inputs and generates relevant recommendations, enhancing the user experience.