**Mini Project Report**                                            Date:12/05/2024

Name: Shashank Rajendra Barai

Matriculation No: 11038167

My journey of developing the chat bot involved a systematic approach after Literature Review of various NLP models, starting with the development of a basic version and then iteratively enhancing its functionality and efficiency through different methods. Initially, I built a simple Chat Bot 1 capable of learning from user input and updating its knowledge base. Then explored text embedding and similarity calculation using a pre-trained Word2Vec model to further refine the chat bot's ability to understand and respond to user queries while continuously learning and improving its knowledge base(Chat Bot 2). Then, finally I progressed using pdf file by integrating Hugging Face's Transformers library and Gradio for a more advanced version, leveraging pre-trained models and optimizing performance through quantization techniques(Chat Bot 3).

**Data Collection and Preprocessing:**

For the first two methods, I gathered information from various sources, including university websites, emails, and FAQ documents. Upon collecting these sources, I compiled the text into a single file. Subsequently, I undertook a cleaning process to extract and consolidate relevant information, ensuring accuracy and coherence. Once the text was refined, I formatted it into a structured JSON file using Python code, facilitating its utilization in the development of chatbots employing JSON-based knowledge bases.

For the third method, I followed a similar process for data acquisition, extracting relevant text from university websites, emails, and FAQ documents. I then compiled and cleaned the text to eliminate irrelevant content and ensure data consistency. Following the cleaning process, I formatted the refined text into a PDF document, which served as the knowledge base for the Transformer-based chatbot.

# Chat Bot 1:

**Preprocessing:**

This script constitutes a basic chatbot that interacts with users by responding to their questions based on a predefined knowledge base stored in a JSON file. The text cleaning or preprocessing steps employed here are minimal. User input is compared against existing questions in the knowledge base using the get_close_matches function from the difflib module to find the best match. This process aids in handling variations in user input, such as misspellings or slight differences in phrasing. However, no explicit text preprocessing steps like stemming, lemmatization, or stop-word removal are performed. When the bot encounters a question it doesn't recognize, it prompts the user to provide

an answer and adds this new question-answer pair to the knowledge base. The only data manipulation involves appending new entries to the existing JSON structure and saving it back to the file.

**Development Process and Solutions Implemented**

**Overview**:

This report outlines the development process and solutions implemented for creating a simple chat bot capable of learning from user input and updating its knowledge base accordingly.

**Development Process:**

1. Requirement Analysis:

- Identified the need for a chat bot capable of answering user queries based on a pre-defined knowledge base.

- Determined the necessity for the bot to dynamically learn from user input and incorporate new information into its knowledge base.

2. **Design Phase:**

- Designed the structure of the chat bot, including functions for loading and saving the knowledge base, finding the best match for user input, and interacting with users.

- Decided to use JSON format for storing the knowledge base due to its simplicity and compatibility with Python.

3. **Implementation:**

- Developed the chat bot in Python, leveraging standard libraries such as json for working with JSON data and difflib for finding close matches to user queries.

- Implemented functions to load and save the knowledge base, find the best match for user input, and interact with users in a loop until the user decides to quit.

4. **Testing and Refinement:**

- Conducted extensive testing to ensure the chat bot functions correctly under various scenarios, including matching user queries accurately and updating the knowledge base appropriately.

- Refinement of the code was done to enhance efficiency, readability, and user experience.

**Solutions Implemented:**

1. **Modular Development:**

- Developed the chat bot using modular design principles, allowing for better organization, reusability, and maintainability of code.

- Encapsulated functionality into separate functions such as loading/saving the knowledge base, finding the best match, and interacting with users.

2. **Dynamic Learning:**

- Implemented a feature for the chat bot to dynamically learn from user input and update its knowledge base with new question-answer pairs.

- Employs the difflib library to find the best match for user queries, facilitating the learning process by identifying similar questions in the knowledge base.

3. JSON-Based Knowledge Base:

- Utilized JSON format for storing the knowledge base, providing a structured and easily accessible storage mechanism.

- Enables seamless loading and saving of the knowledge base, ensuring persistence across multiple sessions of the chat bot.

4. **User Interaction:**

- Designed the chat bot to interact with users in a conversational manner, responding to queries with relevant answers retrieved from the knowledge base.

- Incorporated a loop to facilitate continuous interaction with users until they decide to quit, enhancing the user experience.

**Performance Evaluation:**

In evaluating the effectiveness and accuracy of the chatbot developed through this first method, which involved a structured development process focusing on requirement analysis, design, implementation, testing, and refinement, quantitative metrics played a significant role. I measured the chatbot's response accuracy by calculating the percentage of correctly answered user queries against the total number of queries. Additionally, i assessed the speed and efficiency of the chatbot's responses to ensure optimal performance. Qualitative evaluation involved gathering my bathmates feedback , focusing on aspects such as user satisfaction, perceived accuracy, and overall experience. Through this process, i identified areas for improvement and iteratively refined the chatbot in the next development phase.

**Conclusion:**

The development process resulted in the creation of a functional chat bot capable of answering user queries based on a pre-defined knowledge base while dynamically learning from user input. By employing modular development practices and leveraging appropriate libraries, the chat bot achieves efficiency, accuracy, and user-friendliness.

# Chat Bot 2:

**Preprocessing:**

In this Word2Vec-based chatbot, text preprocessing is streamlined to convert user input and knowledge base questions into dense vector representations using pre-trained word embeddings. Initially, both user input and questions are tokenized by splitting them into words, leveraging Python's `split()` method. These tokens are then individually looked up in the Word2Vec model to obtain their word embedding vectors, forming the basis for semantic understanding. For each text, the word embeddings are averaged to create a single vector representation, facilitating the capture of overall semantic meaning. If any token is absent in the Word2Vec model's vocabulary, a zero vector is utilized to maintain consistent dimensions. This preprocessing workflow ensures efficient similarity calculations between user queries and existing knowledge base questions, enabling accurate responses from the chatbot.

**Development Process and Solutions Implemented**

**Introduction:**

The development process of building a chatbot involved several key steps including acquiring a pre-trained Word2Vec model, designing the chatbot architecture, implementing text embedding and similarity calculation, creating functions to interact with the user, and handling the knowledge base.

1. **Acquiring Pre-trained Word2Vec Model:**

- Utilized the gensim.downloader module to download the pre-trained Word2Vec model trained on the Google News dataset.

- This model was chosen due to its ability to represent words as high-dimensional vectors, capturing semantic relationships between words.

2. **Designing Chatbot Architecture:**

- Developed a modular architecture for the chatbot, comprising functions for loading/saving the knowledge base, embedding text, calculating similarity, finding the best match, and interacting with the user.

3. **Implementing Text Embedding and Similarity Calculation:**

- Implemented the embed_text() function to convert input text into vector representations using the Word2Vec model.

- Developed the calculate_similarity() function to compute cosine similarity between two vectors, facilitating comparison between user input and questions in the knowledge base.

4. **Creating User Interaction Functions:**

- Designed functions to load and save the knowledge base as JSON files.

- Developed functions to find the best match for user input within the knowledge base, and retrieve corresponding answers.

5. **Handling the Knowledge Base:**

- Implemented functionality for the chatbot to learn from user input by asking for answers to questions it doesn't know.

- Updated the knowledge base with new questions and answers provided by the user.

6. **Testing and Refinement:**

- Conducted extensive testing to ensure the chatbot's functionality and accuracy in responding to user queries.

- Refactored code for improved readability, efficiency, and maintainability.

**Performance Evaluation:**

In evaluating the effectiveness and accuracy of the chatbot developed through the 2nd method, which involved acquiring a pre-trained Word2Vec model and implementing text embedding and similarity calculation, our evaluation strategy focused on quantitative metrics and comparative analysis. Quantitatively, I measured response accuracy, and response time to assess the chatbot's performance. Additionally, I conducted comparative analysis across the different development

methods to identify which approach resulted in the most effective and user-friendly chatbot. By evaluating the chatbot's performance against these metrics and comparing it to alternative methods, I gained insights into its effectiveness and accuracy in responding to user queries and moved further for 3rd development phase using transformers

**Conclusion:**

The development process involved acquiring a pre-trained Word2Vec model, designing a modular chatbot architecture, implementing text embedding and similarity calculation, creating functions for user interaction, and handling the knowledge base. Through iterative testing and refinement, the chatbot demonstrates the ability to effectively respond to user queries while continuously learning from user input to enhance its knowledge base.

# Chat bot3:

**Preprocessing**

- Tokenization: The input text is tokenized using the Hugging Face Transformers tokenizer. This step involves breaking down the text into smaller units such as words, subwords, or characters, depending on the tokenizer used.

- Chunking: The input text is chunked into smaller segments. This is done using a recursive character-based text splitter, which divides the text into chunks of a specified size with an overlap. This step is particularly useful when dealing with large documents like PDFs.

- Embedding: The text chunks are embedded into dense vector representations. This is achieved using the Hugging Face embeddings, which encode each chunk into a vector space. These embeddings capture the semantic meaning of the text, enabling downstream tasks like retrieval and generation.

- Vectorization: The embeddings are stored in a vector store. In this script, the Chroma vector store is used, which efficiently indexes and retrieves embeddings based on similarity. This step facilitates fast retrieval of relevant chunks during conversational interactions.

Overall, these preprocessing steps prepare the input text for downstream tasks such as conversational retrieval and text generation, ensuring efficient handling and meaningful interactions in the chatbot system.

**Document Report: Development Process and Implemented Solutions**

1. **Introduction:**

The objective of this document is to provide a comprehensive overview of the development process and solutions implemented for creating a conversational chat bot using Hugging Face's Transformers library and Gradio for user interaction.

2. **Development Process:**

- Package Installation: Necessary packages were installed using pip, including libraries for language processing, UI creation, and model loading.
- Login to Hugging Face Hub: Authentication to the Hugging Face Hub was performed to access pre-trained models.
- Library Import: Required libraries for language processing, model loading, UI creation, and data handling were imported.
- Model Configuration: The conversational model was configured with quantization enabled for efficient inference.
- Prompt Generation: Functions were defined to generate prompts combining system instructions and user instructions.
- Language Processing Setup: Language processing libraries were configured, and PDF documents were loaded for text processing.
- Pipeline Creation: A text generation pipeline was created using the loaded model and tokenizer.
- Class Definition for Chat Bot: A class was defined to handle the chat bot's functionality, including memory, prompt, and retrieval.
- Chat Bot Creation: The chat bot object was created and initialized with a conversational retrieval chain.
- User Interface Setup: Gradio UI components were defined for user interaction, including textboxes, chatbot interface, and buttons.
- Function Definitions: Functions were defined to respond to user inputs, clear memory, and update system prompts.
- Launching the UI: The UI was launched using Gradio with debugging enabled.

3. **Implemented Solutions:**

- Pre-trained Model Usage: Leveraged Hugging Face's Transformers library to access and use pre-trained conversational models, enabling efficient chat bot development.
- Quantization for Efficiency: Implemented quantization techniques to optimize model inference, reducing memory footprint and improving performance.
- Language Processing: Utilized language processing libraries such as MarkdownHeaderTextSplitter and RecursiveCharacterTextSplitter for efficient text processing and document loading.
- User Interface Creation: Developed a user-friendly interface using Gradio, allowing users to interact with the chat bot seamlessly.

- Memory Management: Implemented memory functionality to store conversation history, enabling the chat bot to provide contextually relevant responses.
- Dynamic Prompt Generation: Generated prompts dynamically based on system instructions and user input, enhancing the conversational experience.

**Performance Evaluation:**

For this third development method, which involved leveraging pre-trained models and libraries for efficient chatbot development, my evaluation approach encompassed both quantitative and qualitative measures. Quantitatively, we assessed response accuracy, and response time to gauge the chatbot's performance. I analyzed the rate at which the chatbot learned from user input and updated its knowledge base, providing insights into its adaptability and effectiveness.

4. **Conclusion:**

- The development process involved the integration of various libraries and techniques to create an efficient and user-friendly conversational chat bot.
- Implemented solutions such as pre-trained model usage, quantization for efficiency, language processing, UI creation, and memory management have contributed to the robustness and effectiveness of the chat bot.
- The final product provides an intuitive interface for users to interact with the chat bot, delivering contextually relevant responses based on the conversation history and system instructions.

**What could have been done**:

While I could have opted for gpt based Transformers model using OpenAI's paid API which is fast and ready to make chatbot, I chose to work with pre-trained Word2Vec models for text embedding and similarity calculation and open source language models, which provided me with a learning opportunity in text preprocessing. My decision allowed me to delve deeper into the theoretical aspects of natural language processing while gaining valuable insights into chatbot development. Overall, my choice was driven by a desire for hands-on learning and a focus on foundational concepts, which ultimately enriched my understanding of the subject matter.

-------------------------------------------------------Thank you-----------------------------------------------------------