# Multilingual Speech Recognition Model for RAG

## Solution Approach

**Objective:**

The objective of this project is to develop a multilingual speech recognition and query-Answering system leveraging pre-trained models, specifically focusing on the integration of a multilingual Automatic Speech Recognition (ASR) model with the Retrieval-Augmented Generation (RAG) model. The system aims to automatically detect the language of an audio input, transcribe the speech, translate it to English if necessary, and generate contextually relevant Answers using the RAG model. This approach eliminates the need for additional training of the RAG model for multiple languages, making it a versatile and efficient solution for multilingual applications.

**Background:**

In this project, I utilized a range of advanced tools and technologies to develop a comprehensive multilingual speech recognition system.

- o Retrieval-Augmented Generation (RAG) Model: I employed the RAG model, which combines a retriever and a generator to produce contextually accurate answers. The retriever fetches relevant information from a corpus, while the generator synthesizes responses based on this information, ensuring high-quality and relevant answers.
- o Whisper for Automatic Speech Recognition (ASR):For transcribing spoken language into text, I used OpenAI's Whisper model. This pre-trained model supports multiple languages and accents, providing accurate and reliable transcription of audio inputs.
- o Google Translate for Translation: To handle translations, I integrated the Google Translate API. This tool allows for translating text from the detected language into a target language (e.g., English), enabling effective processing of multilingual inputs.

These technologies were combined to create a system that can accurately transcribe, detect language, translate, and generate responses based on audio queries in various languages.

## 2. System Overview

- **Architecture Diagram:** Include a diagram illustrating the flow from audio input to the final answer generation. This helps visualize how different components interact.

- **Components:** Briefly describe each component of the system:

# Multilingual Speech Recognition Model for RAG

## Solution Approach

| | |
|---|---|
| • Speech-to-Text (ASR) Model | ○ Translation |
| • Language Detection | ○ RAG Model for Answer Generation |

**Solution Approach:**

To accomplish the objective of building a multilingual speech recognition and query-Answering system, I followed a systematic approach leveraging pre-trained models and APIs. First, I initialized the Retrieval-Augmented Generation (RAG) model and its components, including the tokenizer, retriever, and model, ensuring they Ire correctly loaded from the "facebook/rag-sequence-nq" checkpoint. I then set up the Google Translator API and the `langdetect` library for language detection and translation tasks. For speech-to-text conversion, I utilized the pre-trained Whisper model to transcribe audio files, followed by language detection to determine if translation was necessary. If the detected language was not English, I used the Google Translate API to translate the transcribed text into English. The translated text was then tokenized using the RAG tokenizer, and the RAG model generated Answers based on this input. I decoded the outputs to obtain the final Answers. To facilitate user interaction, I created a Gradio interface with inputs for audio files and target language selection, defining the processing function to handle the entire workflow seamlessly. Finally, I launched the Gradio interface, providing users with an intuitive platform to upload audio queries and receive contextually relevant Answers in multiple languages.

**Step 1:** Initialize the RAG Model and Retriever

➢ Load the RAG Model: I utilized Hugging Face's model hub to load the pre-trained RAG (Retrieval-Augmented Generation) model, which combines a retriever with a generator to produce contextually relevant Answers.

➢ Load the Tokenizer: To convert text into a suitable format for the RAG model, I loaded the pre-trained tokenizer associated with the RAG model.

➢ Initialize the Retriever: The retriever component, crucial for fetching relevant documents or context from a dataset, was initialized with appropriate settings, including using a dummy dataset for simplicity.

➢ Error Handling: I ensured proper error handling during the initialization process to catch and report any issues with loading the model, tokenizer, or retriever components.

# Multilingual Speech Recognition Model for RAG

## Solution Approach

**Step 2:** Initialize Google Translator

- ➢ Import the Translator: I utilized the `googletrans` library, which provides an API for translating text between different languages.

- ➢ Create a Translator Instance: An instance of the Translator class was initialized to facilitate text translation.

- ➢ Set Up Language Support: I ensured that the translator could handle a variety of languages, particularly those expected to be encountered in the audio inputs.

- ➢ Error Handling: I implemented error handling to manage potential issues with the translation process, ensuring robustness and reliability.

**Step 3:** Speech-to-Text with Language Detection

- ➢ Check Audio File Existence: I verified that the provided audio file existed to prevent errors during processing.

- ➢ Load the ASR Model: I used a pre-trained Automatic Speech Recognition (ASR) model, Whisper, to convert the audio input into text.

- ➢ Transcribe Audio: The ASR model was applied to the audio file to obtain the transcription of the spoken content.

- ➢ Detect Language: I used a language detection library, `langdetect`, to identify the language of the transcribed text.

- ➢ Error Handling: I ensured robust error handling for issues related to file existence, transcription accuracy, and language detection reliability.

**Step 4:** Translate Text Using Google Translate API

- ➢ Translate Text: I used the Google Translate API to translate the transcribed text from its detected language to the target language (e.g., English).

- ➢ Specify Source and Target Languages: I clearly defined the source language (detected language) and the target language for translation.

- ➢ Handle Translation Response: The translated text was retrieved and processed from the API response.

- ➢ Error Handling: I implemented error handling to manage potential translation errors, including network issues or unsupported languages.

# Multilingual Speech Recognition Model for RAG

## Solution Approach

**Step 5:** Generate Answer with RAG

- ➢ Prepare Input for RAG: The translated text was tokenized using the RAG tokenizer to convert it into a suitable format for the RAG model.
- ➢ Generate Answer: The RAG model was used to generate an Answer based on the tokenized input.
- ➢ Decode the Output: The model's output tokens Ire converted back into human-readable text using the tokenizer.
- ➢ Error Handling: I ensured error handling for potential issues during tokenization, generation, and decoding processes.

**Step 6:** Process Audio Query

- ➢ Detect Language and Transcribe Audio: I utilized the speech-to-text with language detection function to obtain the transcription and detected language from the audio file.
- ➢ Translate Transcription: If the detected language was not English, the transcribed text was translated to English using the translation function.
- ➢ Generate Answer with RAG: The translated text was used as input to the RAG model to generate a relevant Answer.
- ➢ Format the Output: I compiled the detected language, original transcription, translated transcription, and RAG-generated Answer into a formatted response for the user.

**Testing and Evaluation**

To ensure the functionality and effectiveness of the multilingual speech recognition system, I conducted various test cases and evaluated performance metrics. For speech recognition, I tested the system with clear English audio, non-English audio with varying accents, and audio with background noise to validate transcription accuracy. Language detection was assessed using transcribed text in different languages, including less common and mixed languages, to confirm precise language identification. I evaluated translation accuracy by translating text from Spanish and other complex languages, including idiomatic expressions, to ensure high-quality and contextually appropriate translations. Answer generation with the RAG model was tested with both English and translated queries to ensure the relevance and coherence of responses.

# Multilingual Speech Recognition Model for RAG

## Solution Approach

I measured performance using metrics such as Word Error Rate (WER) for speech recognition, detection accuracy and error rates for language identification, BLEU scores and human evaluations for translation quality, and relevance scores and response times for answer generation. These tests and metrics verified that the system operated effectively across all functionalities.

**Conclusion:**

From this project, I can conclude that integrating a pre-trained multilingual Automatic Speech Recognition (ASR) model, such as Whisper, with the Retrieval-Augmented Generation (RAG) model, enables effective multilingual speech recognition and question-Answering capabilities. By implementing speech-to-text conversion, language detection, translation, and context-aware Answer generation, I have created a robust system capable of handling audio inputs in multiple languages. This system successfully transcribes spoken content, translates it if necessary, and generates accurate Answers based on the transcribed text, demonstrating the potential for multilingual applications in various domains.