

# End-to-End AI Voice Assistance

## Introduction

The Speech-to-Speech AI Assistant project aims to create an end-to-end pipeline capable of understanding spoken queries, processing them, and responding with synthesized speech. This assistant is designed to be highly interactive, allowing users to engage with it naturally through voice commands. The system leverages advanced speech recognition, natural language processing (NLP), and text-to-speech (TTS) technologies to achieve seamless communication.

The pipeline's workflow starts with voice activity detection (VAD), which determines when a user is speaking. The detected voice input is then converted to text using automatic speech recognition (ASR). A large language model (LLM) processes the transcribed text to generate an appropriate response. Finally, the system converts this response back into speech using TTS technology, delivering an audible answer to the user.

The project utilizes a variety of state-of-the-art models and libraries, including models for ASR and LLMs such as Mistral and NVIDIA-based solutions, as well as custom modules for voice synthesis, ensuring high performance and flexibility.

## Pipeline Architecture

The AI assistant pipeline consists of the following major components:

1. **Voice Activity Detection (VAD):** Detects when the user starts and stops speaking.
2. **Speech-to-Text (STT):** Converts the user's spoken input into text using a pre-trained model.
3. **Language Model (LLM):** Processes the text input and generates a response using a pre-trained language model.
4. **Text-to-Speech (TTS):** Converts the generated text response back into speech for output.

## Libraries and Tools Used

The project leverages various state-of-the-art machine learning models and tools to implement a voice-to-voice AI assistant. The major components include:

1. **Speech-to-Text (ASR):** Converts spoken input into text using Google's Speech Recognition API.

2. **Language Model (LLM):** Processes the transcribed text to generate a coherent response using the google/gemma-2-2b-it model.
3. **Text-to-Speech (TTS):** Converts the generated response back into speech using parler-tts/parler-tts-large-v1 for natural-sounding voice synthesis with custom voice parameters.

## ASR Implementation

```
1 import speech_recognition as sr
2 import sounddevice as sd
3 import scipy.io.wavfile as wav
4
5 # Recording parameters
6 DURATION = 5 # seconds
7 SAMPLE_RATE = 16000
8
9 def record_audio(filename):
10     print("Recording...")
11     audio = sd.rec(int(DURATION * SAMPLE_RATE), samplerate=SAMPLE_RATE, channels=1, dtype='int16')
12     sd.wait()
13     print("Recording finished.")
14     wav.write(filename, SAMPLE_RATE, audio)
15     print(f"Audio saved as {filename}")
16
17 def transcribe_audio(filename):
18     recognizer = sr.Recognizer()
19     with sr.AudioFile(filename) as source:
20         audio = recognizer.record(source)
21         print("Transcribing...")
22         try:
23             text = recognizer.recognize_google(audio)
24             print("Transcription complete.")
25             return text
26         except sr.UnknownValueError:
27             return "Google Speech Recognition could not understand audio"
28         except sr.RequestError as e:
29             return f"Could not request results from Google Speech Recognition service; {e}"
30
31 def main():
32     audio_file = 'recorded_audio.wav'
33     record_audio(audio_file)
34     transcription = transcribe_audio(audio_file)
35     print(f"Transcribed text: {transcription}")
36
37 if __name__ == "__main__":
38     main()
```

## LLM Response Generation

```
1 import torch
2 from transformers import pipeline
3
4 pipe = pipeline(
5     "text-generation",
6     model="google/gemma-2-2b-it",
7     model_kwargs={"torch_dtype": torch.bfloat16},
8     device="cuda",
9 )
10
11 messages = [
12     {"role": "user", "content": "what is the capital of India."},
13 ]
14
15 outputs = pipe(messages, max_new_tokens=256)
16 assistant_response = outputs[0]["generated_text"][-1]["content"].strip()
17 print([assistant_response])
18
```

## TTS Implementation (using Parler TTS)

```
1 import torch
2 from parler_tts import ParlerTTSForConditionalGeneration
3 from transformers import AutoTokenizer
4 import soundfile as sf
5
6 device = "cuda:0" if torch.cuda.is_available() else "cpu"
7
8 model = ParlerTTSForConditionalGeneration.from_pretrained("parler-tts/parler-tts-large-v1").to(device)
9 tokenizer = AutoTokenizer.from_pretrained("parler-tts/parler-tts-large-v1")
10
11 prompt = "Hey, how are you doing today?"
12 description = "A female speaker delivers a slightly expressive and animated speech with a moderate speed and pitch."
13
14 input_ids = tokenizer(description, return_tensors="pt").input_ids.to(device)
15 prompt_input_ids = tokenizer(prompt, return_tensors="pt").input_ids.to(device)
16
17 generation = model.generate(input_ids=input_ids, prompt_input_ids=prompt_input_ids)
18 audio_arr = generation.cpu().numpy().squeeze()
19 sf.write("parler_tts_out.wav", audio_arr, model.config.sampling_rate)
20
```

## Test audio



recorded\_audio.wav

```
(venv) D:\AIVA2>python -u "d:\AIVA2\record_and_transcribe.py"
Recording...
Recording finished.
Audio saved as recorded_audio.wav
Transcribing...
Transcription complete.
Transcribed text: what is the capital of France
```