

Speech to Text Converter

**A PROJECT REPORT
for
Introduction To AI (AI101B)
Session (2024-25)**

Submitted by

**Abhishek Mishra
202410116100008
Avnish Kaushik
202410116100045
Ankit Kumar Shahi
202410116100028
Bhaskar divedi
202410116100048**

**Submitted in partial fulfilment of the
Requirements for the Degree of**

MASTER OF COMPUTER APPLICATION

**Under the Supervision of
Apoorv Jain
Assistant Professor**



Submitted to

**DEPARTMENT OF COMPUTER APPLICATIONS
KIET Group of Institutions, Ghaziabad
Uttar Pradesh-201206**

(APRIL-2025)

CERTIFICATE

Certified that **Abhishek Mishra 202410116100008, Avnish Kaushik 202410116100045, Ankit Kumar Shahi 202410116100028, Bhaskar divedi 202410116100048** have carried out the project work having “**Speech to Text Converter**” (INTRODUCTION TO AI) (AI101B)) for **Master of Computer Application** from Dr. A.P.J. Abdul Kalam Technical University (AKTU) (formerly UPTU), Lucknow under my supervision. The project report embodies original work, and studies are carried out by the student himself and the contents of the project report do not form the basis for the award of any other degree to the candidate or to anybody else from this or any other University/Institution.

This is to certify that the above statement made by the candidate is correct to the best of my knowledge.

Date: 22-04-2025

Apoorv Jain
Assistant Professor
Department of Computer Applications
KIET Group of Institutions, Ghaziabad
An Autonomous Institutions

Dr. Akash Rajak
Dean
Department of Computer Applications
KIET Group of Institutions, Ghaziabad
An Autonomous Institutions

ABSTRACT

Speech Recognition, also known as Automatic Speech Recognition (ASR), is a transformative technology that enables computers to understand and transcribe human speech into text. This project explores the implementation of a basic speech recognition system using Python and relevant libraries such as speech_recognition, pyaudio, and wave. The system captures audio input through a microphone, processes it, and converts spoken words into readable text with the help of cloud-based and offline APIs.

The project highlights the key components of a speech recognition pipeline including audio capture, voice activity detection, and real-time transcription. The application is capable of recognizing short sentences, commands, or phrases, making it suitable for use cases such as voice-controlled applications, transcription services, and accessibility tools for the differently abled.

Recent advancements in deep learning and cloud computing have significantly improved the performance of speech recognition systems. The simplicity of Python, along with its powerful libraries, provides an ideal platform for building and experimenting with such systems. This report presents the methodology, tools, implementation, and outcomes of building a working ASR model, laying the foundation for more complex and intelligent voice-enabled applications in the future.

ACKNOWLEDGEMENTS

Success in life is never attained single-handedly. My deepest gratitude goes to my project supervisor, Apoorv Jain for his/ her guidance, help, and encouragement throughout my project work. Their enlightening ideas, comments, and suggestions.

Words are not enough to express my gratitude to Dr. Akash Rajak, Professor and Dean, Department of Computer Applications, for his insightful comments and administrative help on various occasions.

Fortunately, I have many understanding friends, who have helped me a lot on many critical conditions.

Finally, my sincere thanks go to my family members and all those who have directly and indirectly provided me with moral support and other kind of help. Without their support, completion of this work would not have been possible in time. They keep my life filled with enjoyment and happiness.

Abhishek Mishra

Avnish Kaushik

Ankit Kumar Shahi

TABLE OF CONTENTS

| S. No. | Content | Page No. |
|--------|--|--------------|
| 1 | Certificate | ii |
| 2 | Abstract | iii |
| 3 | Acknowledgements | iv |
| 4 | Chapter 1: Introduction | 1-2 |
| | 1.1 Overview | 1 |
| | 1.2 Problem Statement | 1 |
| | 1.3 Objective | 1 |
| | 1.4 Scope | 1 |
| | 1.5 Features | 1 |
| | 1.6 Hardware and Software Used | 2 |
| 5 | Chapter 2: Literature Review | 3-5 |
| | 2.1 Introduction to Speech Recognition | 3 |
| | 2.2 Historical Background | 3 |
| | 2.3 Recent Advancements | 4 |
| | 2.4 Technologies Used in Literature | 4 |
| | 2.5 Applications in Previous Works | 5 |
| 6 | Chapter 3: Tools and Technologies Used | 6-8 |
| | 3.1 Python Programming Language | 6 |
| | 3.2 SpeechRecognition Library | 6 |
| | 3.3 PyAudio | 7 |
| | 3.4 Google Speech Recognition API | 7 |
| | 3.5 Additional Python Modules | 8 |
| | 3.6 Hardware Requirements | 8 |
| | 3.7 Software Requirements | 8 |
| 7 | Chapter 4: System Design and Architecture | 9-10 |
| | 4.1 System Overview | 9 |
| | 4.2 Input Module | 9 |
| | 4.3 Silence Detection Module | 9 |
| | 4.4 Speech Recognition Module | 10 |
| 8 | Chapter 5: Implementation Details | 11-12 |
| | 5.1 Introduction | 11 |
| | 5.2 Importing Required Libraries | 11 |
| | 5.4 Audio Stream Initialization | 11 |
| | 5.5 Continuous Listening | 11 |
| | 5.6 Real-Time Output Display | 12 |
| | 5.7 Save Audio to File (Optional) | 12 |

| S. No. | Content | Page No. |
|-------------------|---|-----------------|
| 9 | Chapter 6: Testing and Evaluation | 13-14 |
| | 6.1 Introduction | 13 |
| | 6.2 Testing Environment | 13 |
| | 6.3 Test Scenarios | 14 |
| | 6.4 Evaluation Criteria | 14 |
| | 6.5 Sample Test Result | 14 |
| | 6.6 Observations | 14 |
| | 6.7 Limitations Identified | 14 |
| 10 | Chapter 7: Applications and Use Cases | 15-16 |
| | 7.1 Accessibility and Assistive Technology | 15 |
| | 7.2 Virtual Assistants and Smart Devices | 15 |
| | 7.3 Transcription Services | 15 |
| | 7.4 Call Center Analytics | 15 |
| | 7.5 Language Learning and Training | 15 |
| | 7.6 Medical and Legal Transcription | 16 |
| 11 | Chapter 8: Conclusion and Future Scope | 17-18 |
| | 8.1 Conclusion | 17 |
| | 8.2 Future Scope | 18 |
| 12 | Chapter 9: References | 19 |

Chapter 1

Introduction

1.1 Overview

Speech is the most natural and efficient form of communication for humans. With rapid advancements in artificial intelligence and natural language processing, speech recognition systems have become increasingly prevalent in real-world applications. A Speech to Text Converter is a software application that captures spoken words through a microphone and converts them into written text using speech recognition algorithms. This project leverages Python and various powerful libraries like `speech_recognition`, `pyaudio`, and `wave` to build a simple yet effective speech-to-text system that can process real-time voice input.

1.2 Problem Statement

Typing or writing text manually is often time-consuming and can be challenging for people with disabilities or those multitasking in dynamic environments. Traditional input methods also limit hands-free interaction with devices. This project addresses the need for a more intuitive, accessible, and efficient form of text input through voice commands.

1.3 Objective

The primary objective of this project is to design and implement a Python-based application that:

- Captures real-time speech using a microphone,
- Processes the audio data using speech recognition APIs,
- Converts the spoken content into accurate textual format,
- Supports both real-time and recorded audio inputs.

1.4 Scope

This project is limited to converting English speech to text using local microphones and pre-recorded .wav files. The system supports basic audio processing, silence detection, and works efficiently in environments with moderate noise levels. It does not support speaker identification or advanced linguistic understanding.

1.5 Features

- Real-time voice input via microphone
- Conversion of .wav audio files to text
- Silence detection and noise adjustment
- Compatibility with multiple speech APIs (Google, IBM, Sphinx)

1.6 Hardware and Software Used in Project

- **Hardware:**
 - Microphone
 - Computer with minimum 4GB RAM
- **Software:**
 - Python 3.x
 - Libraries: `speech_recognition`, `pyaudio`, `wave`, `audioop`, `os`, `time`
 - Operating System: Windows/Linux

Chapter 2

Literature Review

2.1 Introduction to Speech Recognition

Speech Recognition, also known as Automatic Speech Recognition (ASR), is a transformative technology that allows machines and computers to interpret and convert spoken language into written text. This technology is designed to bridge the gap between human speech and machine understanding, enabling computers to process natural language input in real time. The process typically involves capturing the audio signals of spoken words, analyzing the patterns in the sound waves, and using algorithms to identify the corresponding words or phrases. This field of study has garnered considerable attention for its potential applications in various domains such as transcription services, virtual assistants, voice command systems, and even healthcare.

Recent advancements have drastically improved the accuracy, speed, and versatility of ASR systems. With continuous research in natural language processing (NLP), neural networks, and big data, modern speech recognition systems are able to handle complex tasks such as real-time transcription, multilingual support, and even understanding of non-native accents. These systems are now being used in an array of applications, from virtual assistants like Siri, Google Assistant, and Alexa, to automated customer service systems, live captioning, and healthcare documentation. As a result, ASR technology has become an integral part of the digital transformation, enabling more intuitive and efficient interactions between humans and machines.

Despite its impressive progress, challenges still remain in the field, including issues with accuracy in noisy environments, recognition of multiple speakers, and contextual understanding. However, ongoing innovations continue to push the boundaries of what ASR systems can achieve, promising even greater accuracy, accessibility, and user experience in the future.

2.2 Historical Background

The journey of speech recognition began in the 1950s at Bell Laboratories, where early pioneers in the field worked to develop systems capable of understanding spoken language. One of the first notable projects was the Audrey system, which could recognize spoken digits, specifically the numbers from 0 to 9, spoken by a single speaker. This was a groundbreaking achievement, as it marked one of the first successful attempts to teach machines to recognize human speech.

The 1960s and 1970s witnessed further efforts to expand the capabilities of speech recognition systems. One of the most significant contributions during this period came from the Defense Advanced Research Projects Agency (DARPA), which launched the Speech Understanding Research (SUR) program in the early 1970s. This project, which ran from 1971 to 1976, was instrumental in advancing speech recognition beyond simple isolated words to more complex forms of speech. DARPA's SUR program brought together top academic and research institutions, including Carnegie Mellon University, SRI International, and BBN Technologies, to collaborate on developing new technologies for speech understanding.

The advancements made during the SUR program were crucial in pushing the boundaries of speech recognition. They set the stage for the next generation of systems capable of handling larger vocabularies, speaker-independent recognition, and even continuous, natural speech. These early efforts

by DARPA and other research organizations paved the way for the more advanced ASR systems we use today, marking a pivotal moment in the evolution of speech recognition technology.

2.3 Recent Advancements

Modern speech recognition has been significantly revolutionized by the advent of Deep Neural Networks (DNNs), which allow systems to learn intricate patterns in acoustic signals and make more accurate predictions. Recurrent Neural Networks (RNNs), particularly Long Short-Term Memory (LSTM) networks, further enhanced the technology by addressing the challenge of capturing long-range dependencies in speech, allowing for better understanding of context over time. In more recent years, Transformer-based models, such as OpenAI's Whisper and Google's Speech-to-Text API, have taken speech recognition to new heights, leveraging self-attention mechanisms to process entire sequences of speech simultaneously, improving both speed and accuracy. These cutting-edge models have shown great potential in handling noisy environments and multiple speakers. Additionally, cloud-based solutions now provide real-time transcription and multilingual support, enabling seamless integration into applications across different languages without the need for local processing power, making speech recognition more accessible and scalable than ever before.

2.4 Technologies Used in Literature

Many speech recognition systems rely on the following tools and frameworks:

- **SpeechRecognition (Python Library):** Acts as a unified façade over multiple ASR engines—Google Cloud Speech-to-Text, IBM Watson, CMU Sphinx, Microsoft Azure, and more—so you can swap providers without rewriting your audio-handling code. It manages microphone buffering, WAV-file loading, and API key configuration, and exposes a simple `recognize_<engine>()` interface along with built-in retry logic for network hiccups or “no speech detected” errors.
- **PyAudio:** A thin Python wrapper around the native PortAudio library that provides real-time, low-latency access to your system’s audio devices. You can configure sample rates (e.g., 16 kHz for speech), channel counts (mono vs. stereo), and frame-buffer sizes to balance CPU usage against responsiveness—crucial when you need to capture unbroken streams of spoken input.
- **wave & audioop (Standard Libraries):**
 - The `wave` module lets you read, write, and inspect WAV files—handy for saving raw microphone captures for later analysis or debugging.
 - The `audioop` module offers byte-level signal processing routines (RMS, sample-width conversion, bias removal), which are the building blocks of silence detection, volume normalization, and even simple effects like trimming leading/trailing silence before sending audio to an ASR engine.
- **Voice Activity Detection (VAD):** An algorithmic “gatekeeper” that examines short windows of audio—often via RMS energy or spectral entropy—to decide if someone is speaking. By dropping silent or background-only frames, VAD reduces wasted API calls,

cuts transcription costs, and keeps your application responsive by only processing chunks that actually contain voice.

2.5 Applications Discussed in Previous Works

- **Accessibility: Helping Users with Disabilities Interact Using Voice:** Speech recognition plays a vital role in making technology accessible for individuals with disabilities. For people with motor impairments, voice commands allow them to interact with devices without using traditional input methods like keyboards or touchscreens. Similarly, people with visual impairments can benefit from voice-controlled navigation and screen-reading technologies, enabling them to access information and perform tasks more independently.
- **Smart Assistants: Siri, Google Assistant, Alexa:** Speech recognition forms the backbone of personal digital assistants like Apple's Siri, Google's Google Assistant, and Amazon's Alexa. These assistants rely on accurate voice recognition to understand user queries and provide real-time responses. By using natural language processing (NLP) and machine learning, smart assistants not only respond to simple commands but can also handle complex tasks like setting reminders, controlling smart home devices, playing music, or even providing traffic updates.
- **Call Centers: Automatic Call Routing and Response Analysis:** In call centers, speech recognition helps automate call routing by identifying the purpose of the call or extracting key information from the conversation. For instance, it can detect keywords like "billing," "support," or "sales," directing the call to the appropriate department. Furthermore, speech recognition technology is increasingly used for response analysis—analyzing customer-agent conversations to assess service quality, detect sentiment, and generate insights for better decision-making.
- **Transcription Services:** Converting Lectures, Interviews, and Podcasts to Text: Transcription services powered by speech recognition have made it easier to convert spoken content into written form, saving time and effort compared to manual transcription. This technology is widely used for converting lectures, interviews, podcasts, and meetings into accurate text. It is especially helpful in education and research, where transcribing long hours of audio or video content can be a time-consuming task. Furthermore, the increasing accuracy of speech-to-text engines ensures that transcriptions are reliable, enabling content to be searchable and more accessible.

The literature reveals an increasing reliance on cloud APIs and machine learning-based solutions. While older systems were limited to specific vocabularies and environments, newer systems are versatile, scalable, and language-independent.

Chapter 3:

Tools and Technologies Used

3.1 Programming Language: Python

Python is a high-level, interpreted language known for its readability and vast ecosystem of libraries. For this project, Python was chosen due to its simplicity and powerful libraries for speech processing, such as:

- **speech_recognition:** The speech_recognition library is a powerful Python module that facilitates easy access to multiple speech-to-text APIs. This library allows you to interface with popular speech recognition services like Google Web Speech API, Microsoft Bing Voice Recognition, IBM Speech to Text, and even offline engines like CMU Sphinx. It abstracts much of the complexity, enabling developers to focus on integration rather than underlying technical details. This makes it a go-to choice for projects involving automatic speech recognition (ASR).
- **pyaudio:** PyAudio is a cross-platform library that allows Python programs to interact with audio hardware. It provides simple and efficient methods to capture real-time audio from microphones and output sound through speakers. For this project, PyAudio is used to stream microphone input to the system, which is then processed by the speech recognition engine. The ability to handle continuous audio streams with low latency makes PyAudio ideal for real-time speech recognition applications.
- **wave:** The wave module is a built-in Python library that simplifies working with WAV files, a common format for storing raw audio. It provides functions for reading and writing audio data in the WAV format, making it easy to save and process audio recordings from the microphone. For this project, the wave library is used to handle audio input before it's passed to the speech recognition engine. It can also be used to export transcribed audio for further analysis or storage.
- **os and time (for system operations and delays):** The os and time libraries provide essential system-level functionality that allows the program to interact with the operating system. The os module is used for file management tasks such as checking if an audio file exists, creating directories, or setting environmental variables. On the other hand, the time module is crucial for introducing delays, timing operations, or managing timeouts during the speech recognition process. For example, it can be used to introduce pauses between repeated attempts to capture speech or to set the maximum length of time to wait for input before stopping.

3.2 SpeechRecognition Library

This is the core library used in the project. It supports multiple APIs and engines for speech recognition. The main features include:

- **Easy Integration with Google's Web Speech API:** The speech_recognition library allows seamless integration with Google's Web Speech API, enabling fast and accurate speech-to-text conversion with minimal setup.
- **Support for Noise Adjustment:** It offers built-in noise filtering and pre-processing to handle background noise, improving recognition accuracy in noisy environments.
- **Microphone Input Support:** The library supports real-time microphone input via PyAudio, making it easy to capture audio directly from the microphone for live speech recognition.
- **Compatibility with WAV Files:** It works well with WAV audio files, allowing easy reading, writing, and manipulation of audio data in the WAV format for later processing or storage.

3.3 PyAudio

PyAudio is a cross-platform audio I/O library used to record real-time audio from the microphone. In this project, it helps:

- **Stream Audio Directly into the System:** PyAudio enables the direct streaming of audio input from the microphone into the system, allowing for real-time audio capture and processing without delays.
- **Configure Sample Rates, Channels, and Chunk Sizes:** PyAudio allows customization of audio parameters such as sample rate (e.g., 16 kHz for speech), channels (mono or stereo), and chunk sizes (the number of frames per buffer), giving you control over audio quality and performance.
- **Manage Audio Buffers in Real-Time:** PyAudio handles real-time audio buffers efficiently, enabling the continuous capture and processing of audio data without interruptions or memory issues. This ensures smooth operation for live speech recognition applications.

3.4 Google Speech Recognition API

This cloud-based API processes spoken words and returns transcriptions. It supports over 100 languages and has features like:

- **Automatic Punctuation:** The API automatically inserts commas, periods, and question marks in transcribed text, producing more readable output without manual editing.
- **Noise Filtering:** Built-in noise reduction algorithms attenuate background sounds, ensuring the transcription focuses on the speaker's voice for higher accuracy.
- **Real-Time and Recorded Input Processing:** Supports both live streaming audio (for immediate transcription) and prerecorded files (for batch processing), allowing flexible integration into diverse workflows.

Note: Internet access is required for this API to work.

3.5 Additional Python Modules

- **wave:** The wave module is used to read, write, and manipulate .wav audio files in Python. It allows easy access to raw audio data for processing or saving audio recordings captured from a microphone.
- **audioop:** The audioop module is useful for performing audio-related operations such as measuring the RMS (Root Mean Square) of audio chunks. RMS is often used for detecting silence or low-volume segments, enabling more accurate speech recognition by focusing on meaningful speech parts.
- **os & time:** The os module handles file operations such as checking for the existence of files, creating directories, or managing paths, while the time module is used for introducing delays or controlling time-related operations. Together, they enable efficient system-level tasks and synchronization in the speech recognition workflow.

3.6 Hardware Requirements

- **Microphone:** To capture audio input
- **Computer:** With Python environment installed (preferably with at least 4 GB RAM)
- **Stable Internet Connection:** For accessing the Google Speech API

3.7 Software Requirements

- **Operating System:** Windows 10/Linux/macOS
- **Python Version:** 3.7 or above
- **IDEs Used:** VS Code / Jupyter Notebook
- **Required Libraries:** Installed via pip install SpeechRecognition pyaudio wave

Chapter 4

System Design and Architecture

4.1 System Overview

The Real-Time Speech Recognition System is designed to continuously listen to audio input from the microphone, detect silence to optimize transcription, convert speech to text using the Google API, and display the output in real time. The system architecture is modular and scalable for future enhancements like multilingual support or offline recognition.

4.2 Input Module

- **Device: Microphone** The microphone serves as the primary input device, capturing real-time audio from the user. It is typically a condenser or dynamic microphone, chosen for its clarity and sensitivity in capturing speech. The microphone converts sound into electrical signals, which are then converted into digital data by the audio interface or sound card.
- **Function: Captures Real-Time Audio Data** The microphone continuously records audio as long as the input stream is active. This data is then converted into digital samples in real-time, ready for processing by the system for tasks like speech recognition. The captured audio is streamed directly into the software without any delays, allowing for seamless interaction.
- **Library: PyAudio** : The PyAudio library is used to interface with the microphone and handle the streaming of audio data. It provides simple methods to initialize and manage the input stream, configure audio parameters (such as sample rate and channels), and capture live audio samples for further processing. PyAudio is efficient and easy to use, making it an excellent choice for real-time applications.
- **Features:** Set to Mono-Channel for Speech Clarity and Minimal Processing To streamline the process and minimize computational complexity, the microphone input is configured to mono-channel. This ensures that all audio data is captured in a single track, which simplifies the speech recognition process and reduces the need for additional processing required for stereo input. Mono-channel audio also ensures clearer speech capture, as it avoids unnecessary complexity and noise associated with multiple channels.

4.3 Silence Detection Module

- **Technique:** Root Mean Square (RMS) amplitude is used to determine silence.
- **Library:** audioop
- **Function:** The system records only when the RMS value is above a threshold to reduce unnecessary processing and API usage.

4.4 Speech Recognition Module

- **API Used:** Google Speech Recognition API via speech_recognition library.
- **Process:**
 - Audio is passed as chunks to the recognizer.
 - API returns transcribed text.
 - Errors like UnknownValueError and RequestError are handled gracefully.

Chapter 5

Implementation Details

5.1 Introduction

The implementation of the Real-Time Speech Recognition System was done using Python, leveraging its powerful libraries for audio input and speech processing. This chapter provides a detailed explanation of the code and its functionality.

The project contains a single Python file named speech_recognition_live.py with clearly defined blocks for initialization, recording, recognition, and output display.

5.2 Importing Required Libraries

```
import speech_recognition as sr  
import os  
import time  
import audioop  
import wave  
import pyaudio
```

- speech_recognition: Handles interaction with the Google Speech Recognition API.
- audioop: For silence detection using RMS values.
- wave: To save audio if needed.
- pyaudio: Captures real-time microphone input.

5.4 Audio Stream Initialization

```
r = sr.Recognizer()  
mic = sr.Microphone()  
  
• Recognizer() is the engine that handles all recognition tasks.  
• Microphone() accesses the system's mic.
```

5.5 Continuous Listening with Silence Detection

with mic as source:

```
r.adjust_for_ambient_noise(source)
```

```
print("Listening...")
while True:
    audio = r.listen(source, phrase_time_limit=5)
    try:
        text = r.recognize_google(audio)
        print(f"You said: {text}")
    except sr.UnknownValueError:
        print("Could not understand audio")
    except sr.RequestError as e:
        print(f"Request failed; {e}")
```

- Adjusts for background noise.
- Listens in chunks of 5 seconds.
- Sends audio to Google API.
- Handles exceptions.

5.6 Real-Time Output Display

The transcribed text is displayed in the console immediately after being processed.

Example Output:

```
yaml
CopyEdit
Listening...
You said: Hello, how are you?
You said: This is a speech recognition test.
```

5.7 Optional: Save Audio to File (for debugging or training data)

```
wf = wave.open("output.wav", 'wb')
wf.setnchannels(1)
wf.setsampwidth(p.get_sample_size(pyaudio.paInt16))
wf.setframerate(44100)
wf.writeframes(b''.join(frames))
wf.close()
```

This section is currently commented but can be enabled to save audio input as a .wav file.

Chapter 6

Testing and Evaluation

6.1 Introduction

Testing is an essential part of any software development process. For the Real-Time Speech Recognition System, both functional and performance tests were carried out to ensure the system operates as expected under various conditions.

6.2 Testing Environment

- **Device Used:** HP 14s Laptop
- **Processor:** Intel i3 (10th Gen)
- **RAM:** 8 GB
- **Operating System:** Windows 10
- **Internet Connection:** 50 Mbps
- **Python Version:** 3.10
- **Microphone:** In-built system mic

6.3 Test Scenarios

| Test Case ID | Description | Expected Result | Actual Result | Status |
|--------------|-----------------------------------|------------------------------------|--------------------|--------|
| TC01 | Speak a clear sentence in English | Transcription matches speech | Accurate | ✓ |
| TC02 | Speak with background noise | Partial or correct transcription | Minor words missed | ✓ |
| TC03 | Remain silent | System should wait or skip | No text output | ✓ |
| TC04 | Speak non-English words | Return unclear or incorrect output | As expected | ✓ |
| TC05 | Disconnect Internet | Raise RequestError | Handled gracefully | ✓ |
| TC06 | Continuous 1-minute speech | Maintain real-time performance | Smooth output | ✓ |

6.4 Evaluation Criteria

The system was evaluated on:

- **Accuracy:** Compared transcribed text with actual speech.
- **Speed:** Observed delay between speech and output.
- **Robustness:** Tested with multiple edge cases (noise, silence, long speech).
- **Usability:** Checked whether beginners could easily use the system.

6.5 Sample Test Result

```
In [6]: # Start Listening
listen_and_transcribe()

Listening... (Press Ctrl+C to stop)

Speak now...
>>> hello what are you doing

Speak now...
(Could not understand audio)

Speak now...

Stopped listening
```

Speech Input:

“Hello What are you doing.”

Transcription Output:

“Hello What are you doing.”

Accuracy: 100%

6.6 Observations

- The system performed best in quiet environments.
- Recognition accuracy decreased slightly in noisy areas.
- Delay was minimal (~1-2 seconds), acceptable for most real-time uses.
- The system worked well even with mid-range hardware and basic microphones.

6.7 Limitations Identified

- Requires stable internet for API requests.
- Does not support offline recognition (can be added via vosk or pocketsphinx).
- Accuracy drops for accents or fast speech.

Chapter 7

Applications and Use Cases

7.1 Accessibility and Assistive Technology

- **Voice-to-Text for the Hearing Impaired:** Real-time transcription empowers people with hearing loss to “read” conversations as they happen in classrooms, meetings, or one-on-one discussions.
- **Hands-Free Control for Mobility-Impaired Users:** Users who cannot operate keyboards can issue voice commands to navigate software, type emails, or control smart home devices.

7.2 Virtual Assistants and Smart Devices

- **Integration with Assistants (e.g., Alexa, Google Home):** The same principles—microphone input, silence detection, cloud-based ASR—underlie consumer voice assistants that control lights, play music, or fetch information.
- **IoT and Home Automation:** Embedding a real-time recognition module in an IoT hub allows users to turn on/off appliances, adjust thermostats, and manage security systems via simple voice phrases.

7.3 Transcription Services

- **Meeting and Lecture Capture:** Automatic transcription of boardroom meetings or university lectures saves time on manual note-taking and creates searchable text archives.
- **Media and Journalism:** Reporters can use the system for rapid dictation of briefs, interviews, and live reporting with minimal lag.

7.4 Call Center Analytics

- **Real-Time Agent Assistance:** As customer service representatives speak with callers, the system transcribes dialogue and can trigger knowledge-base searches or recommended responses on the agent’s screen.
- **Quality Monitoring and Sentiment Analysis:** Transcriptions feed NLP pipelines to extract key metrics—call duration, customer sentiment, compliance with scripts—for supervisors and analytics teams.

7.5 Language Learning and Pronunciation Training

- **Immediate Feedback for Learners:** Language apps can record a learner’s pronunciation, transcribe it, and compare against a model sentence to highlight pronunciation errors.

- **Accent Adaptation Research:** Collecting large volumes of learner speech and transcriptions helps build models that adapt ASR to non-native accents, thus improving global usability.

7.6 Medical and Legal Transcription

- **Doctor's Dictations:** Physicians dictate patient notes during or after examinations; real-time transcription can populate electronic health records instantly.
- **Courtroom Proceedings:** Transcribing witness testimonies and judge's remarks in real time reduces the lag in producing official transcripts.

Chapter 8

Conclusion & Future Scope

8.1 Conclusion

The Real-Time Speech Recognition System presented in this project successfully demonstrates the ability to convert speech into text with minimal delay and high accuracy under most conditions. By leveraging existing cloud-based ASR technologies, such as the Google Speech-to-Text API, the system provides a robust solution for various real-time transcription needs.

Key achievements of the project include:

- **Accurate Transcription:** The system consistently delivered accurate transcriptions, even with minor background noise.
- **Real-Time Processing:** The application worked effectively for short bursts of speech, with only minimal delay (~1–2 seconds).
- **Scalability:** The modularity of the system allows for the easy addition of features such as multi-language support or an offline recognition engine.

The system provides a strong foundation for future development, and its potential applications span a wide range of industries, including accessibility for the hearing impaired, real-time transcription in professional environments, and voice-activated smart devices.

8.2 Future Scope

The current system, while functional, has several areas for improvement and future enhancements:

1. **Multi-Language Support:**
 - The current version only supports English. Implementing additional language models or integrating with multilingual APIs will extend its usability to non-English-speaking regions and international markets.
2. **Offline Speech Recognition:**
 - As the system heavily relies on cloud-based services, its performance is dependent on internet connectivity. Integrating an offline ASR system, like Vosk or PocketSphinx, would provide more flexibility and reduce dependency on an internet connection.
3. **Enhanced User Interface (UI):**
 - Currently, the system only outputs transcriptions in a console-based format. Developing a GUI (Graphical User Interface) would make it more user-friendly and accessible to a broader audience.
4. **Advanced Noise Filtering:**
 - The system's performance can degrade in noisy environments. Implementing more advanced noise cancellation techniques or custom models trained on diverse speech data could significantly improve accuracy in such scenarios.

5. Support for Continuous Speech:

- Long-form speech, such as speeches or meetings, can cause delays or inaccuracies. Future versions could include features like automatic segmentation of continuous speech into smaller chunks or real-time punctuation insertion to improve the readability of the transcriptions.

6. Custom Vocabulary and Domain-Specific Models:

- To further improve accuracy, especially for technical fields like medical or legal transcription, custom vocabulary sets and domain-specific models could be trained and integrated into the system.

7. Integration with Other Applications:

- The system can be integrated with other software, such as word processors, content management systems, or customer service platforms, for automatic documentation generation or real-time support.

Chapter 9

References

1. Rabiner, L. R., & Juang, B. H. (1993). *Fundamentals of Speech Recognition*. Prentice-Hall, Inc.
2. Jurafsky, D., & Martin, J. H. (2023). *Speech and Language Processing* (3rd ed.). Draft available at: <https://web.stanford.edu/~jurafsky/slp3/>
3. OpenAI. (2022). *Whisper: Robust Speech Recognition via Foundation Models*. Retrieved from: <https://openai.com/research/whisper>
4. Google Cloud. (n.d.). *Cloud Speech-to-Text Documentation*. Retrieved from: <https://cloud.google.com/speech-to-text>
5. Python Software Foundation. (n.d.). *SpeechRecognition Library Documentation*. Retrieved from: <https://pypi.org/project/SpeechRecognition/>
6. PyAudio Documentation. (n.d.). Retrieved from: <https://people.csail.mit.edu/hubert/pyaudio/>
7. Wikipedia Contributors. (2024). *Speech Recognition*. In *Wikipedia, The Free Encyclopedia*. Retrieved from: https://en.wikipedia.org/wiki/Speech_recognition
8. Mozilla. (n.d.). *DeepSpeech Documentation*. Retrieved from: <https://github.com/mozilla/DeepSpeech>
9. IBM. (n.d.). *Watson Speech to Text*. Retrieved from: <https://www.ibm.com/cloud/watson-speech-to-text>
10. Microsoft Azure. (n.d.). *Speech Service Documentation*. Retrieved from: <https://learn.microsoft.com/en-us/azure/cognitive-services/speech-service/>