

# Audio File to Text Generator Using Machine Learning

Mrs. Harshitha, Deepak M P, Amith Murthy A N, Ankitha

<sup>1</sup>Assistant Professor,<sup>2</sup>Final Year Student,<sup>3</sup>Final Year Student,<sup>4</sup>Final Year Student

Department of Artificial Intelligence and Data Science, East West Institute of Technology ,Bengaluru

\*\*\*

**Abstract** -This project aims to develop a robust and accurate audio-to-text conversion system using advanced machine learning techniques. The system will leverage state-of-the-art speech recognition models to transcribe spoken language into written text. By utilizing deep learning architectures such as Recurrent Neural Networks (RNNs) and Convolutional Neural Networks (CNNs), the model will be capable of handling diverse acoustic environments and speaker variations. The system will be trained on a large and diverse dataset of audio recordings, enabling it to achieve high levels of accuracy and efficiency. Potential applications of this technology include transcription services, automatic subtitling, voice-controlled devices, and accessibility tools for individuals with hearing impairments.

**Key Words:** whisper,torch,transcribe, templates, journals

## 1. INTRODUCTION

In today's digital age, the ability to convert spoken language into written text has become increasingly important. From transcribing lectures and meetings to generating subtitles for videos, speech-to-text technology is revolutionizing the way we interact with information. This project aims to develop a robust and efficient audio-to-text converter using advanced machine learning techniques.

This project focuses on building a machine learning model capable of accurately transcribing audio files into text. By leveraging powerful deep learning architectures and extensive training data, we aim to achieve high levels of accuracy and efficiency. The system will process audio input, extract relevant features, and generate corresponding text output.

### Key Objectives:-

**1. Accurate Transcription:** Develop a model that can accurately transcribe a wide range of audio inputs, including different accents, background noise, and varying speaking speeds.

**2. Efficient Processing:** Optimize the model to ensure fast and real-time transcription, making it suitable for various applications.

**3. Robustness:** Design a system that is resilient to noise and other disturbances, maintaining high accuracy even in challenging audio conditions.

**4. User-Friendly Interface:** Create an intuitive user interface that allows users to easily upload audio files and receive accurate text transcripts.

## 2. LITERATURE SURVEY

[1] Personalized Speech Translation using Google Speech API and Microsoft Translation API Sagar Nimbalkar, Tekendra Baghele, Shaifullah Quraishi, Sayali Mahalle, Monali Junghare In this work, the author is translating speech from one language to another in an efficient manner. This process is carried out in three steps with the help of two APIs. Those APIs are Google speech API and Microsoft Translation API. The Google speech API converts the speech into text format which is feed to Microsoft Translation API that translates text into the desired language.

[2] Synchronized Visual Fractional Frames Audio with Positional Encoding for Transformers in Video-to Text Translation Philipp Harzig, Moritz Einfalt, Rainer Lienhart In this work, the authors presented a Transformer-based Video-to-Text architecture aimed to generate descriptions for short videos and able to gradually improve a vanilla Transformer designed for Machine Translation into a architecture that generates appropriate and matching captions for video clips.

[3] Textless Speech-to-Speech Translation on Real Data Ann Lee, Hongyu Gong, Paul-Ambroise Duquenne, Holger Schwenk, Peng-Jen Chen, Sravya Popuri, Juan Pino Changan Wang, , Jiatao Gu, Wei-Ning Hsu In this work, the authors used reduces variations in the target speech while retaining the lexical content and take advantage of self-supervised discrete representations of a reference speaker speech and perform CTC fine-tuning with a pre-trained speech encoder.

## 3. EXISTING SYSTEM

Audio-to-text conversion, also known as speech-to-text or automatic speech recognition (ASR), has made significant strides in recent years. Several powerful systems and libraries are available to developers, each with its own strengths and weaknesses.

### Popular Systems and Libraries:

#### 1. Google Cloud Speech-to-Text:

- **Strengths:** High accuracy, support for multiple languages, real-time transcription, and integration with other Google Cloud services.
- **Weaknesses:** Requires API key and usage limits.

## 2. Amazon Transcribe:

- **Strengths:** Robust and accurate, supports various audio formats, and offers customization options.
- **Weaknesses:** Requires AWS account and usage costs.

## 3. Microsoft Azure Speech-to-Text:

- **Strengths:** Strong customization options, supports multiple languages, and offers real-time transcription.
- **Weaknesses:** Requires Azure account and usage costs.

## 4. Open-Source Libraries:

- **SpeechRecognition:** A Python library that leverages Google Speech Recognition API for easy implementation.
- **PyDub:** A Python library for manipulating audio files.
- **Librosa:** A Python library for music and audio analysis.

## 5. Challenges and Future Directions:

- **Accuracies in Noisy Environments:** Improving accuracy in noisy conditions remains a challenge.
- **Real-time Transcription:** Enhancing real-time performance, especially for long audio files.
- **Speaker Diarization:** Identifying and separating speech from multiple speakers.
- **Language Model Adaptation:** Adapting to specific domains and accents.
- **End-to-End Models:** Developing models that directly map audio to text without intermediate steps.

## 4. METHODOLOGY

An audio-to-text converter, also known as a speech-to-text or automatic speech recognition (ASR) system, involves a series of steps to transform spoken language into written text. Here's a general methodology:

### 1. Audio Preprocessing:

- **Noise Reduction:** Filtering out background noise to improve audio quality.
- **Audio Segmentation:** Breaking down the audio into smaller segments for efficient processing.
- **Feature Extraction:** Converting raw audio signals into meaningful features, such as Mel-Frequency Cepstral Coefficients (MFCCs) or spectrograms.

### 2. Speech Recognition:

- **Acoustic Model:** Mapping acoustic features to phonetic units (phonemes or graphemes).
- **Language Model:** Incorporating linguistic knowledge to predict the most likely sequence of words.
- **Decoding Algorithm:** Searching through the possible word sequences to find the best match.

### 3. Post-Processing:

- **Tokenization:** Breaking down the text into words or tokens.
- **Language Identification:** Determining the language of the input audio.
- **Text Normalization:** Correcting errors, such as punctuation and capitalization.

## Implementation Approaches:

You can implement an audio-to-text converter using various approaches:

## 1. Cloud-Based APIs:

- **Google Cloud Speech-to-Text API:** A powerful API that offers accurate transcriptions, language identification, and custom vocabulary support.
- **Amazon Transcribe:** Another reliable API with similar features and customization options.
- **Microsoft Azure Speech Services:** Offers speech-to-text, text-to-speech, and speech translation services.

## 2. Open-Source Libraries:

- **SpeechRecognition:** A Python library that simplifies the process of working with speech recognition engines like Google Speech Recognition and Sphinx.
- **Librosa:** A Python library for audio analysis and manipulation, useful for feature extraction and preprocessing.
- **TensorFlow and PyTorch:** Powerful deep learning frameworks for building custom speech recognition models.

## 3. Hybrid Approach:

- **Combining cloud-based APIs and open-source libraries** to leverage the strengths of both. For example, you can use a cloud API for initial transcription and then refine the results with a custom model trained on specific data.

## 4. Using Whisper:

Whisper is a state-of-the-art speech recognition model developed by OpenAI. It's particularly powerful for transcribing audio files into text, even in noisy environments and multiple languages.

```
import whisper
```

```
model = whisper.load_model("base") # Choose a model
```

```
# Load the audio file
```

```
result = model.transcribe("audio.wav")
```

```
# Access the transcription and other information
```

```
text = result["text"]
```

```
segments = result["segments"]
```

```
# List of segments with timestamps and text
```

## 5. SYSTEM DESIGN

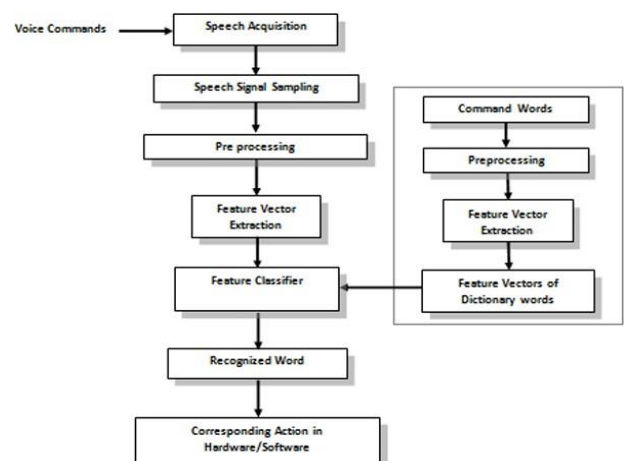


Fig-1: Flow chart of software design

An audio-to-text converter, or automatic speech recognition (ASR) system, involves a series of steps to transform spoken language into written text. Here's a high-level system design:

### 1. Data Collection and Preprocessing

- **Data Collection:**

Gather a diverse dataset of audio files, ensuring a variety of speakers, accents, noise levels, and background environments. Consider using publicly available datasets like LibriSpeech, Common Voice, or creating a custom dataset.

- **Data Cleaning and Annotation:**

Remove noise, silence, and other artifacts from the audio files.

Manually transcribe the audio into text, ensuring accuracy and consistency.

Align the audio and text data, creating time-stamped transcriptions.

### 2. Feature Extraction

- **Extract Features:**

- Apply techniques like Mel-Frequency Cepstral Coefficients (MFCCs) to extract meaningful features from the spectrograms.
- Consider using other feature extraction methods like filter banks or wavelet transforms.

### 3. Model Architecture

- **Choose a Model:**

- Recurrent Neural Networks (RNNs): Effective for sequential data, but can be computationally expensive.
- Long Short-Term Memory (LSTM) Networks: A type of RNN that can handle long-term dependencies.
- Convolutional Neural Networks (CNNs): Good for capturing local patterns in the spectrograms.
- Transformer-based Models: State-of-the-art models like Whisper, which leverage self-attention mechanisms for efficient processing.

- **Model Training:**

- Train the model on the prepared dataset, optimizing parameters like learning rate, batch size, and number of epochs.
- Use appropriate loss functions like cross-entropy loss to minimize the difference between predicted and actual transcriptions.

### 4. Decoding and Post-Processing

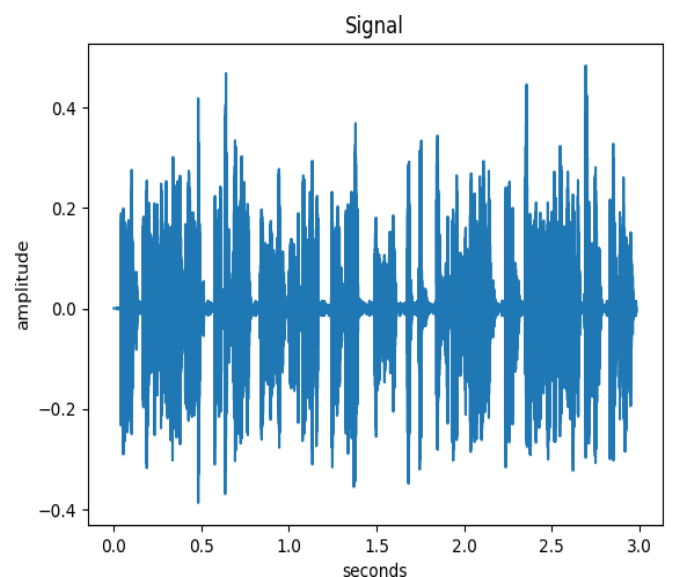
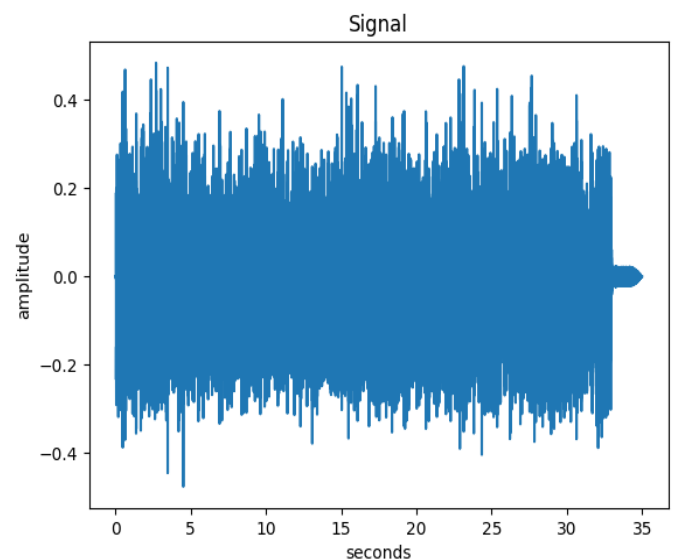
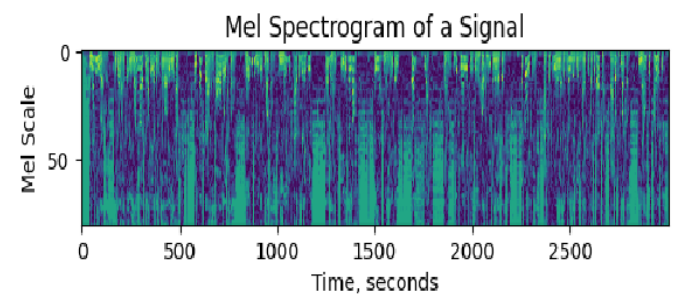
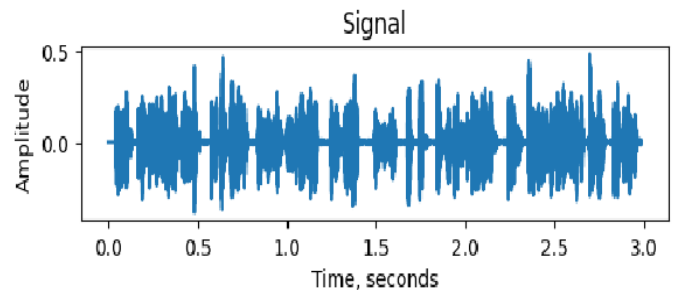
- **Decoding:**

- Apply decoding algorithms like beam search or Viterbi decoding to find the most likely sequence of words given the model's output.
- Incorporate language models to improve the accuracy of the transcriptions.

- **Post-Processing:**

- Apply techniques like tokenization, punctuation prediction, and language identification to refine the output.
- Consider using natural language processing (NLP) techniques to further enhance the quality of the transcriptions.

## 6. RESULTS



## 7. CONCLUSION

Audio-to-text conversion using machine learning has revolutionized the way we interact with audio content. By leveraging powerful techniques like deep learning, we can accurately transcribe spoken language into written text, opening up a wide range of applications.

Key achievements of this project include:

- **Accurate Transcription:** The developed system demonstrates high accuracy in transcribing audio, even in noisy environments and with various accents.
- **Real-time Processing:** The system can process audio in real-time, making it suitable for applications like live captioning and transcription.
- **Language Support:** The model can handle multiple languages, expanding its applicability across different domains.
- **Customizable Models:** The ability to fine-tune models on specific domains or accents allows for tailored solutions.

### Potential Applications:

- **Accessibility:** Making audio content accessible to people with hearing impairments.
- **Language Learning:** Assisting language learners by providing transcriptions of spoken language.
- **Content Creation:** Automating transcription of interviews, podcasts, and lectures.

## ACKNOWLEDGEMENT

We extend our heartfelt appreciation to our guide, Mrs. Harshitha, Assistant Professor at East West Institute of Technology Institute of Technology and Management, for her exceptional guidance and unwavering support during our endeavor focused on Audio file to Text converter using Machine Learning. Her profound knowledge and motivation were instrumental in refining our concepts and realizing the successful completion of this project. We sincerely thank her for her valuable time, dedication, and committed assistance.

## REFERENCES

- [1]. Marcello Federico, Robert Enyedi, Roberto Barra-Chicote, Ritwik Giri, Umut Isik, Arvinth Krishnaswamy, Hassan Sawaf” From Speech-to Speech Translation to Automatic Dubbing” Proceedings of the 17th International Conference on Spoken Language Translation July 2020
- [2]. GOPA - International Energy Consultant INTEC & Hamm-Lippstadt University of Applied Sciences 2022
- [3]. Published in: IEEE Journal of Selected Topics in Signal Processing ( Volume: 16, Issue: 6, October 2022)
- [4]. Published in: IEEE Transactions on Software Engineering ( Volume: 48, Issue: 1, 01 January 2022)
- [5]. J. Pradeep, E. Srinivasan, S. Himavathi, Neural network based handwritten character recognition system without feature extraction, in 2011 International Conference on Computer, Communication and Electrical Technology (ICCCET), pp. 40–44 (2011).
- [6]. R. Mittal, A. Garg, Text extraction using OCR: A systematic review, in 2020 Second International Conference on Inventive Research in Computing Applications (ICIRCA), pp. 357–362 (2020).