# International Journal of Research Publication and Reviews

# Converting Text to Speech and Speech to Text Using Natural Language Processing

## *V. Sai Shashank[1], I. Sweety Julia[2], K. Sai Prasad[3], P. Vennela[4], [5]R. Sai Prasad, K . Sai Rithick[6]*

1,2,3,4,5,6 [CSE]-AI&ML, MRUH

[1]sashank.varanasi@gmail.com, [2]sweety.juliaimmanuel@mallareddyuniversity.ac.in, [3]kavalisaiprasad37@gmail.com

[4]vennelapsb@gmail.com, [5]saiprasadakhil@gmail.com, [6]kommidisairithick@gmail.com

## ABSTRACT:

This project aims to develop an integrated system that seamlessly translates speech to text (STT) and text to speech (TTS) using cutting-edge Natural Language Processing (NLP) techniques. The primary objective is to create a robust system capable of converting spoken language into written text and then back into natural-sounding speech, with a special focus on multilingual capabilities. Once the speech is transcribed, the text will be refined and processed using advanced NLP tools like NLTK or SpaCy. These tools will enhance the quality of the text by addressing grammatical errors, contextual inconsistencies, and semantic nuances, making the text more coherent and contextually appropriate.For the text-to-speech synthesis, the project will employ technologies like Google's Text-to-Speech API. This will ensure that the generated speech sounds natural and human-like, with appropriate intonation and emphasis, thereby enhancing the user experience.The ultimate goal of this project is to significantly improve human-computer interaction by providing a seamless and intuitive interface for users. Applications of this technology are vast, including assistive technologies for individuals with disabilities, language learning tools, and real-time translation services that can bridge communication gaps in multilingual environments.

Keywords: STT, TTS, NLP, NLTK.

## 1.Introduction

The "Speech-to-Text with Sentiment Analysis and Text-to-Speech" project addresses the need for an integrated system that seamlessly converts spoken language into written text, analyzes its sentiment, and provides an audible response. The rise of voice-activated applications and AI-driven customer service tools highlights the growing demand for advanced speech processing solutions. However, most current solutions often lack combined functionality, such as real-time sentiment analysis of spoken input followed by responsive, emotion-aware feedback.

This project aims to fill this gap by creating an application that not only converts speech to text but also determines the sentiment of the speech content (positive, negative, or neutral). The output sentiment is then relayed back to the user using text-to-speech technology, providing a fully interactive and intuitive user experience. The implementation of this project has the potential to support various applications, including customer service automation, personal assistance, and enhanced accessibility tools for individuals with speech or reading difficulties. It brings speech recognition, sentiment analysis, and speech synthesis together in a single, user-friendly interface.

## 2.Literature Review

This paper presents a pioneering approach using deep neural networks (DNNs) for speech recognition tasks, eliminating the need for traditional feature extraction steps. Graves et al. demonstrate how long short-term memory (LSTM) networks improve recognition accuracy for continuous speech, showcasing significant advancements in speech transcription accuracy. Their end-to-end model structure was instrumental in modern STT applications, laying the groundwork for real-time speech recognition systems by simplifying complex processes.[1]

Tacotron is an advanced neural network model for text-to-speech synthesis that converts text directly into audio. Shen et al. introduce an architecture that eliminates traditional speech processing steps, directly mapping text sequences to mel spectrograms. The Tacotron model, later refined in Tacotron 2, generates high-quality synthetic speech that closely mimics human intonation and pronunciation, a breakthrough for achieving natural-sounding TTS systems and enabling applications where responsive and personalized audio output is critical.[2]

This paper introduces the Valence Aware Dictionary for Sentiment Reasoning (VADER), a sentiment analysis model optimized for short, informal text. Hutto and Gilbert designed VADER to handle social media language, making it suitable for conversational analysis. The model incorporates lexicon and

rule-based scoring, considering word intensity and punctuation to achieve a compound sentiment score. VADER's effectiveness in detecting sentiments has made it popular in real-time applications where quick sentiment feedback is essential.[3]

Deng and Fu investigate various approaches to speech emotion recognition (SER), analyzing feature extraction techniques like Mel-frequency cepstral coefficients (MFCC) and prosodic features, combined with algorithms such as SVMs and neural networks. The authors explore how voice parameters, like pitch and tone, correlate with emotions, enabling SER systems to classify emotional states like happiness, anger, and sadness. Their insights are foundational for integrating emotional detection into STT systems, as they highlight methods for capturing sentiment from speech alone.[4]

## 3.Dataset

In this project, the datasets utilized primarily revolve around the functionalities of speech recognition and sentiment analysis. These datasets facilitate the training and functioning of the models that power the application. **Speech Recognition Dataset** While the application leverages Google's Speech Recognition API for converting spoken language to text, the underlying dataset that powers this API consists of a vast collection of audio samples across different languages, accents, and dialects. **Sentiment Analysis Dataset** For sentiment analysis, the VADER sentiment analysis model does not require a specific dataset for initial training since it is a lexicon and rule-based approach. However, it leverages sentiment-labeled datasets for validation and enhancement of its performance.

## 4. Methodology

**Speech Recognition API** and **Google STT** for accurate transcription.

**VADER Sentiment Analysis** to categorize input as positive, negative, or neutral.

**pyttsx3** for synthesizing text into human-like speech

The "Speech-to-Text with Sentiment Analysis and Text-to-Speech" application utilizes a combination of methods and algorithms to achieve its functionalities. Each component employs specific techniques that work together to deliver accurate and efficient results. For the speech-to-text functionality, the application relies on the Google Speech Recognition API. This API uses advanced deep learning models to transcribe spoken language into text. It incorporates techniques like acoustic modeling, language modeling, and phonetic recognition to ensure high accuracy across various accents and languages. The model continuously learns and adapts, improving its performance over time with more data.

In sentiment analysis, the VADER (Valence Aware Dictionary and sEntiment Reasoner) algorithm is utilized. VADER operates on a lexicon-based approach, which assigns sentiment scores to words based on their emotional valence. The algorithm analyzes the context of the words and takes into account modifiers that can alter sentiment intensity, such as negations and degree modifiers. This makes it particularly effective for analyzing social media text and conversational language, where context plays a crucial role. For the text-to-speech conversion, the Pyttsx3 library is employed. This library uses a speech synthesis engine that converts the analyzed text back into audible speech. The synthesis process involves generating speech waveforms from the text input, utilizing pre-recorded phonemes and rules for speech production. The flexibility of Pyttsx3 allows for various voice options and the ability to adjust speech rate and volume.

## 5.Design

The design of the "Speech-to-Text with Sentiment Analysis and Text-to-Speech" application revolves around creating a streamlined user experience that integrates multiple complex processes in a simple and cohesive manner. The design goal is to offer an intuitive, interactive platform where users can effortlessly speak, analyze the emotional content of their speech, and receive audio feedback based on the analysis. The application design includes three core components: the speech-to-text module, the sentiment analysis engine, and the text-to-speech output, all linked through an interactive web interface developed with Streamlit. The architecture leverages a modular design, allowing each component to operate independently while working in conjunction with the others.

### 5.1 Deployment using Stream lit

The code is deployed using streamlit for frontend to provide user interface to convert text to speech and speech to text
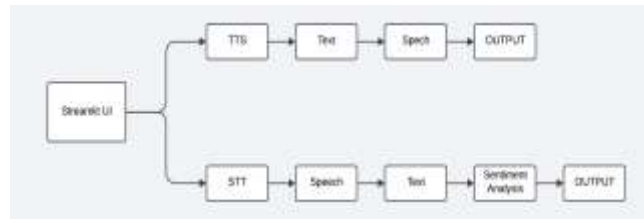
*5.2 Architecture*



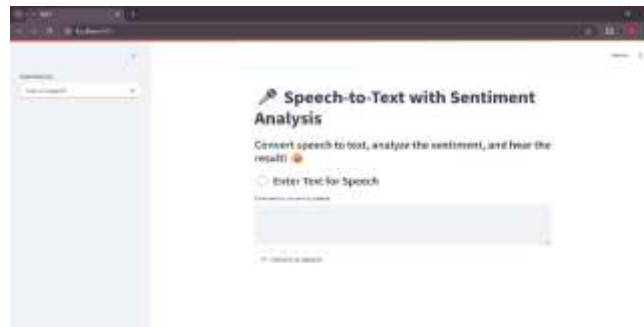Fig 1. Architecture

*5.3 Output Screens*
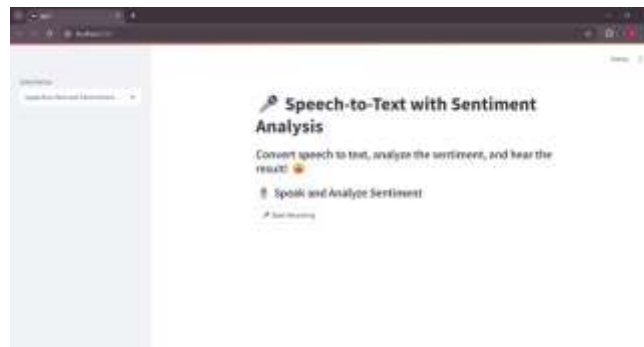


Fig 2. Frontend for TTS



Fig 3 . Frontend for STT
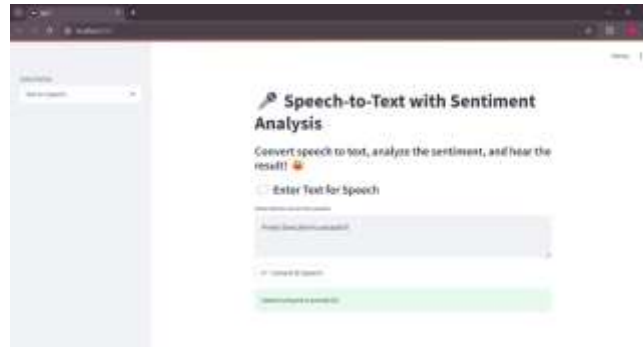


Fig 4.  Converting speech to text

Fig 5. Converting text to speech

## 6. Conclusion

The "Speech-to-Text with Sentiment Analysis and Text-to-Speech" application successfully integrates advanced technologies to provide users with an interactive and efficient tool for communication and analysis. The project has achieved its primary objectives by delivering accurate speech recognition, insightful sentiment analysis, and high-quality text-to-speech conversion. Through the use of state-of-the-art libraries and APIs, the application ensures seamless performance across all functionalities. The implementation of the Google Speech Recognition API allows for real-time transcription with minimal error rates, while the VADER sentiment analysis tool effectively captures the emotional nuances of spoken and written language. Additionally, the integration of Pyttsx3 for text-to-speech conversion enhances user engagement by providing clear and natural voice output.

## 7. Future Work

The "Speech-to-Text with Sentiment Analysis and Text-to-Speech" application has laid a solid foundation for further enhancements and expansions. The future scope of this project encompasses several potential improvements and new features that can enhance its functionality, usability, and accessibility.

**Multilingual Support** One of the most significant areas for expansion is the incorporation of multilingual support. By integrating additional language models into the speech recognition and text-to-speech functionalities, the application can cater to a more diverse user base. This enhancement would allow users to interact with the application in their preferred language, significantly broadening its accessibility and usability. **Improved Sentiment Analysis Techniques** Future iterations of the application could benefit from the integration of more advanced sentiment analysis techniques, including deep learning models. Exploring natural language processing (NLP) frameworks like TensorFlow or PyTorch could enable the development of models that understand contextual sentiments more deeply. This advancement would enhance the accuracy of sentiment classification, particularly in nuanced or complex conversational scenarios. **User Personalization Options** Adding features for user personalization could significantly enhance the user experience. Allowing users to select preferred voice types, speech rates, and even tone settings would provide a more tailored interaction. Additionally, incorporating user profiles to save preferences and interaction history could foster a more engaging experience.

## 8. References:

[1] "End-to-End Speech Recognition Models Using Deep Learning" by Alex Graves, Abdel-rahman Mohamed, Geoffrey Hinton in 2013.

[2] "Tacotron: Towards End-to-End Speech Synthesis" by Jonathan Shen, Ruoming Pang, Ron J. Weiss in 2017.

[3] "VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text" by C.J. Hutto, Eric Gilbert in 2014.

[4] "Speech Emotion Recognition: Features and Algorithms" by Anna W. Deng, Szu-Wei Fu in 2016.

[5] "End-to-End Multilingual Speech Recognition System" by T. Sainath, Yu Zhang, Daniel Rybach in 2018.

[6] "Transfer Learning for Text-to-Speech Synthesis" by Aron Szabadi, Shrikanth Narayanan in 2020.