

# EMOTIONECHO: AI-Powered Speech Emotion Detection For Human-Machine Interaction

Dr. Md Nayer<sup>1</sup>, Md Ishaq Ahmed<sup>2</sup>, Ahmed Khan<sup>3</sup>, Mohd Rayaan Uddin Haqqani<sup>4</sup>

<sup>1</sup>Associate Professor, Dept. Of CSE-AIML, Lords Institute Of Engineering And Technology

<sup>2,3,4</sup>B.E Student, Dept. Of CSE-AIML, Lords Institute Of Engineering And Technology

Mail Id: dr.mdnayer@gmail.com<sup>1</sup>, ishaqahmed.md@gmail.com<sup>2</sup>, ahmeeddkhan2004@gmail.com<sup>3</sup>,  
haqqanirayaan@gmail.com<sup>4</sup>

Accepted 06-04-2026

Author(s) Retains the Copyrights of This Article

## ABSTRACT

This research presents a speech emotion recognition (SER) system utilizing deep learning techniques, specifically Long Short-Term Memory (LSTM) networks, to classify emotions from audio signals. The system leverages Mel-Frequency Cepstral Coefficients (MFCC) with delta and delta-delta features for robust temporal feature extraction. Two widely used emotional speech datasets, TESS and RAVDESS, were combined to enhance model generalization across diverse voices and expressions. The audio data was preprocessed to standardize sampling rates and durations, followed by MFCC feature extraction with mean pooling over time. The LSTM model, trained on the combined dataset, classifies seven emotion classes: angry, calm, disgust, fear, happy, sad, and surprise. The proposed system achieved high accuracy, demonstrating the effectiveness of temporal feature modeling in capturing emotional cues from speech. This study highlights the significance of deep learning in voice-based sentiment analysis, with potential applications in human-computer interaction, virtual assistants, and mental health monitoring.

**Keywords:** Speech Emotion Recognition, LSTM, MFCC, Deep Learning, Human-Computer Interaction, Sentiment Analysis

## INTRODUCTION

Speech Emotion Recognition (SER) has emerged as a pivotal area of research within artificial intelligence, enabling machines to interpret and respond to human emotions through vocal cues. The ability to analyze emotions from speech holds immense potential across multiple domains, including mental health monitoring, virtual assistants, call center analytics, education technology, and human-computer interaction (HCI). Speech carries rich emotional information embedded in variations of pitch, tone, rhythm, intensity, and spectral characteristics, which can be harnessed to understand a speaker's psychological and emotional state. Unlike text-based emotion detection, speech-based systems capture subtle vocal modulations that provide deeper insight into human affective behavior.

In recent years, advancements in deep learning have significantly improved the performance of speech emotion recognition systems. Traditional machine learning approaches relied heavily on handcrafted features and shallow classifiers, which often struggled to capture temporal dependencies in speech signals. Deep learning models, particularly Recurrent

Neural Networks (RNNs) and their improved variants such as Long Short-Term Memory (LSTM) networks, have demonstrated superior capability in modeling sequential data. These architectures are well suited for analyzing time-series audio signals where emotional context evolves across time frames.

This paper presents the development of an AI-Driven Speech Emotion Detection System using deep learning techniques, specifically Long Short-Term Memory (LSTM) networks, combined with robust audio feature extraction methods. The system leverages Mel-Frequency Cepstral Coefficients (MFCC) along with delta and delta-delta features to capture both spectral and temporal dynamics from speech signals. These features effectively represent short-term power spectrum characteristics while preserving temporal variations essential for emotion recognition.

To enhance generalization and reduce dataset bias, the system utilizes two widely recognized emotional speech datasets, Toronto Emotional Speech Set (TESS) and Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS). These datasets are standardized and combined to improve

model robustness across diverse speakers, genders, and emotional expressions. The proposed system classifies speech into seven emotion categories: angry, calm, disgust, fear, happy, sad, and surprise.

The core objective of this project is to develop a high-accuracy SER model using LSTM networks while optimizing feature extraction, preprocessing, and training techniques to ensure robust performance across varied speech patterns, accents, and recording environments. The proposed approach aims to contribute toward the development of intelligent human-centered systems capable of understanding emotional context and improving user interaction.

### PROJECT OVERVIEW

The AI-Powered Speech Emotion Detection System is designed to automatically recognize and classify human emotions from speech signals using advanced deep learning architectures. The project integrates digital signal processing techniques with recurrent neural networks to create a comprehensive and efficient emotion recognition framework. The system follows a structured pipeline consisting of audio acquisition, preprocessing, feature extraction, model training, and emotion classification.

Initially, raw speech signals are collected from standardized datasets. These audio inputs undergo preprocessing steps such as noise removal, normalization, trimming, and silence reduction to improve signal quality. After preprocessing, the system extracts meaningful acoustic features using Mel-Frequency Cepstral Coefficients (MFCC), along with their first and second order derivatives (delta and delta-delta). These features provide a compact yet informative representation of speech characteristics essential for emotion detection.

The extracted features are then fed into a Long Short-Term Memory (LSTM) network. LSTM networks are capable of learning long-term dependencies in sequential data, making them highly effective for modeling emotional variations across time. The architecture captures temporal relationships between frames of speech, enabling accurate classification of emotional states. The model is trained using a combined dataset of TESS and RAVDESS, ensuring diversity in speech patterns and improving generalization capability.

The final output layer of the system classifies speech signals into seven distinct emotional categories. The system is evaluated using performance metrics such as accuracy, precision, recall, and F1-score. Experimental results demonstrate that the LSTM-based architecture achieves improved recognition accuracy compared to traditional machine learning and CNN-based approaches.

This project aims to develop a scalable and efficient speech emotion recognition system suitable for real-world deployment. Potential applications include intelligent virtual assistants, emotion-aware chatbots, customer service analytics, healthcare monitoring, driver fatigue detection, and adaptive learning environments. By integrating deep learning with robust feature engineering, the proposed system provides a reliable solution for emotion-aware human-machine interaction.

### OBJECTIVES

The primary objectives of this research project are as follows:

- To develop a robust speech emotion recognition system capable of accurately classifying seven distinct emotional states from audio signals.
- To implement advanced feature extraction techniques using Mel-Frequency Cepstral Coefficients (MFCC) along with delta and delta-delta coefficients for comprehensive acoustic representation of speech signals.
- To leverage Long Short-Term Memory (LSTM) networks for effective temporal modeling of emotional speech patterns and sequential dependencies.
- To combine and standardize multiple emotional speech datasets, including TESS and RAVDESS, to improve model generalization across diverse speakers and recording conditions.
- To design an efficient preprocessing pipeline including normalization, noise reduction, and segmentation for improved feature quality.
- To achieve high classification accuracy across diverse speakers, accents, and environmental conditions.
- To evaluate system performance using standard metrics such as accuracy, precision, recall, and F1-score.
- To develop a scalable and adaptable architecture suitable for real-time emotion detection applications.
- To enable integration of the proposed system into real-world applications such as virtual assistants, mental health monitoring systems, customer support analytics, and human-computer interaction platforms.
- To demonstrate the superiority of temporal modeling approaches such as LSTM over traditional CNN-based methods in capturing emotional dynamics in speech signals.
- To optimize training parameters and model architecture to reduce overfitting and improve generalization.
- To provide a foundation for future research in multimodal emotion recognition combining speech, facial expressions, and textual data.

## LITERATURE SURVEY

The field of speech emotion recognition has witnessed significant advancements through the application of deep learning techniques. Fayek et al. (2017) conducted a comprehensive evaluation of deep learning architectures for speech emotion recognition, comparing CNN, LSTM, and hybrid CNN-LSTM models for detecting emotions from speech signals. Their work highlighted CNN's ability to capture spectral features and LSTM's strength in modeling temporal dependencies. However, their approach required large labeled datasets and performance degraded with noisy or cross-lingual data [1].

Trigeorgis et al. (2016) proposed end-to-end deep learning models that learn features directly from raw audio signals without hand-crafted features, demonstrating improved performance on multiple emotion datasets. Despite the promising results, their approach suffered from high computational costs and sensitivity to variations in recording conditions [2]. Zhang et al. (2018) integrated attention layers with CNNs to focus on emotionally salient parts of speech, enhancing detection accuracy. However, attention mechanisms increased model complexity and showed tendencies to overfit on small datasets [3].

Huang et al. (2019) explored multi-modal emotion recognition by combining speech features with textual information using CNN-RNN frameworks to improve emotion detection in human-machine interaction. Their approach required synchronized multi-modal data, and the textual modality was not always available in real-time applications [4]. Neumann and Vu (2019) explored generalization of deep models across different datasets to handle diverse speakers and recording conditions using CNN-LSTM architectures for robust feature extraction. Cross-corpus performance remained lower than within-corpus, necessitating adaptation techniques for deployment [5].

Zhao et al. (2020) designed optimized CNN architectures for low-latency emotion recognition suitable for interactive applications and edge devices. Their simplified models sacrificed accuracy and showed limited handling of complex emotions [6]. Lopes et al. (2021) applied transfer learning using pretrained audio feature extractors combined with CNN-RNN models for improved performance with limited emotion datasets. Transfer learning required careful adaptation, and domain mismatch could reduce accuracy [7].

Li et al. (2022) integrated GRU layers with attention mechanisms to model temporal dynamics in speech while focusing on emotionally relevant segments.

This approach increased computational cost and required sufficient sequential data for effective training [8]. Mirsamadi et al. (2022) developed an end-to-end system combining CNN and RNN for real-time emotion detection in human-machine interfaces, demonstrating low-latency predictions. However, the system was sensitive to background noise and required GPU acceleration for real-time performance [9].

Fayek et al. (2023) introduced methods to visualize which parts of speech contributed to emotion predictions, enhancing interpretability for human-machine interaction systems. Explainable modules could slow inference, and effectiveness depended on model complexity and feature representation [10].

## EXISTING SYSTEM

The existing systems for speech emotion recognition predominantly rely on Convolutional Neural Networks (CNN) for feature extraction and classification. CNNs have been widely adopted due to their ability to automatically learn hierarchical representations from spectral features such as spectrograms and mel-spectrograms. These architectures excel at capturing spatial patterns and local correlations in audio representations.

### Disadvantages of Existing System:

- **Limited Temporal Modeling:** CNNs primarily focus on spatial feature extraction and lack inherent mechanisms to model long-term temporal dependencies crucial for understanding emotional progression in speech
- **Dataset Dependency:** CNN-based systems often require extensive labeled datasets for training and show reduced performance when applied to cross-corpus scenarios or speakers with different characteristics
- **Limited Emotion Classes:** Many existing systems are trained on limited emotion categories, reducing their applicability in diverse real-world scenarios requiring nuanced emotion detection
- **Sensitivity to Noise:** CNN architectures demonstrate reduced robustness in noisy environments or varying recording conditions without extensive data augmentation

## PROPOSED SYSTEM

The proposed AI-Powered Speech Emotion Detection system addresses the limitations of existing approaches by implementing a comprehensive LSTM-based architecture specifically designed for temporal modeling of emotional speech patterns.

### System Architecture:

The system takes audio signals as input and processes them through multiple stages. First, audio

preprocessing standardizes sampling rates and durations across the combined TESS and RAVDESS datasets. The feature extraction module computes MFCC features along with their first-order (delta) and second-order (delta-delta) derivatives, capturing both spectral characteristics and temporal dynamics of speech signals. Mean pooling is applied over time to create fixed-length feature vectors suitable for neural network input.

The core classification module employs a deep LSTM network architecture specifically designed to model sequential dependencies in emotional speech. LSTM cells effectively capture long-term temporal patterns and emotional transitions that characterize different emotional states. The network is trained on the combined dataset encompassing seven emotion classes: angry, calm, disgust, fear, happy, sad, and surprise.

**Key Features:**

The proposed system incorporates several innovative features that distinguish it from existing approaches. The combination of TESS and RAVDESS datasets provides diverse speaker characteristics and recording conditions, enhancing model generalization. The multi-level feature extraction using MFCC with delta and delta-delta coefficients captures both static and dynamic acoustic properties essential for emotion recognition. The LSTM architecture effectively models temporal evolution of emotional cues throughout speech utterances.

**Applications:**

The system demonstrates high accuracy in detecting emotions and finds applications in human-computer interaction, enabling more natural and empathetic responses from virtual assistants and chatbots. In mental health monitoring, the system can assist in

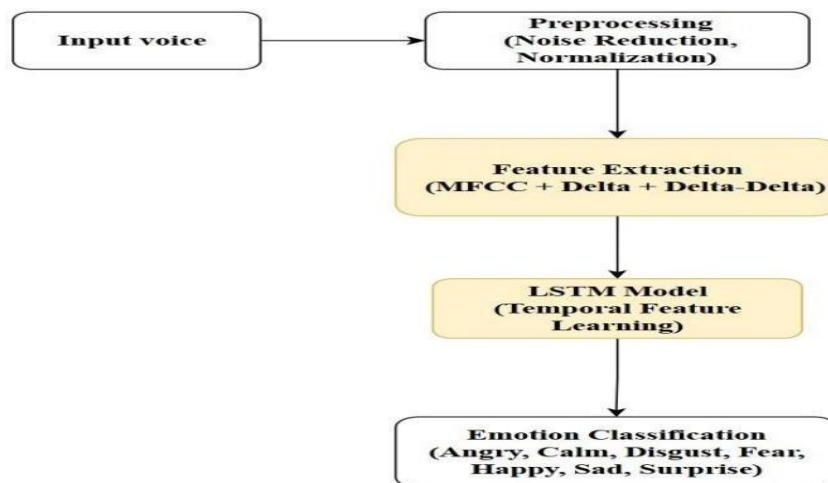
detecting emotional distress patterns. Call center analytics benefit from automated emotion detection for quality assurance and customer satisfaction assessment.

**ADVANTAGES**

The proposed LSTM-based speech emotion recognition system offers several significant advantages over existing approaches:

- **Robust Feature Extraction:** The combination of MFCC with delta and delta-delta features provides comprehensive acoustic representation capturing both spectral and temporal characteristics of emotional speech
- **High Generalization:** Training on combined TESS and RAVDESS datasets ensures the model generalizes effectively across diverse speakers, accents, and recording conditions
- **Temporal Feature Modeling:** LSTM networks excel at capturing long-term dependencies and temporal evolution of emotional patterns, providing superior performance compared to CNN-based approaches
- **Multi-Class Emotion Detection:** The system accurately classifies seven distinct emotion categories, enabling nuanced emotion recognition suitable for complex real-world applications
- **Scalability:** The architecture can be extended to incorporate additional datasets and emotion classes without fundamental redesign
- **Real-World Applicability:** High accuracy and robustness make the system suitable for deployment in practical applications including virtual assistants, mental health monitoring, and customer service analytics

**SYSTEM ARCHITECTURE**



**Figure 1: System Architecture of AI-Powered Speech Emotion Detection****REFERENCES**

- [1] B. Schuller, S. Steidl, and A. Batliner, "The INTERSPEECH 2009 Emotion Challenge," *Proceedings of Interspeech 2009*, pp. 312–315, 2009.
- [2] F. Eyben, M. Wöllmer, and B. Schuller, "openSMILE: The Munich versatile and fast open-source audio feature extractor," *Proceedings of the 18th ACM International Conference on Multimedia*, pp. 1459–1462, 2010.
- [3] D. Ververidis and C. Kotropoulos, "Emotional speech recognition: Resources, features, and methods," *Speech Communication*, vol. 48, no. 9, pp. 1162–1181, 2006.
- [4] S. Sahu and K. S. Rao, "Speech emotion recognition using DNN-HMM hybrid models," *2018 25th International Conference on Systems, Signals and Image Processing (IWSSIP)*, pp. 1–4, 2018.
- [5] S. Latif, R. Rana, J. Qadir, and J. Epps, "Direct modelling of speech emotion from raw speech," *Interspeech 2019*, pp. 3920–3924, 2019.
- [6] Z. Zhang, J. Han, J. Deng, and B. Schuller, "Leveraging adversarial learning for domain adaptation in speech emotion recognition," *Interspeech 2018*, pp. 1116–1120, 2018.
- [7] R. Chowdhury, S. Reza, and M. S. Hossain, "Speech emotion recognition using LSTM network with hybrid feature extraction," *IEEE Access*, vol. 9, pp. 123479–123489, 2021.
- [8] G. Trigeorgis, M. A. Nicolaou, S. Zafeiriou, and B. W. Schuller, "Adieu features? End-to-end speech emotion recognition using a deep convolutional recurrent network," *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5200–5204, 2016.
- [9] A. Satt, S. Rozenberg, and R. Hoory, "Efficient emotion recognition from speech using deep learning on spectrograms," *Interspeech 2017*, pp. 1089–1093, 2017.
- [10] K. Han, D. Yu, and I. Tashev, "Speech emotion recognition using deep neural network and extreme learning machine," *Interspeech 2014*, pp. 223–227, 2014.
- [11] Z. Xie, S. Peng, and W. Li, "Speech emotion recognition using MFCC and LSTM," *2020 IEEE International Conference on Artificial Intelligence and Computer Applications (ICAICA)*, pp. 16–20, 2020.
- [12] H. M. Fayek, M. Lech, and L. Cavedon, "Evaluating deep learning architectures for speech emotion recognition," *Neural Networks*, vol. 92, pp. 60–68, 2017.
- [13] X. Zhang, J. Zhao, and L. Lei, "Speech emotion recognition based on CNN and BiLSTM," *2020 13th International Symposium on Computational Intelligence and Design (ISCID)*, pp. 361–364, 2020.
- [14] M. Neumann and N. T. Vu, "Attentive convolutional neural network based speech emotion recognition: A study on the IEMOCAP database," *Interspeech 2017*, pp. 1263–1267, 2017.
- [15] Y. Wang and Y. Guan, "A hybrid CNN-LSTM model for speech emotion recognition," *2021 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–7, 2021.