# Acoustic Analysis for Early Detection of Vocal Disorders

## Mr.Shiva Kumar[1], Y. Rahul[2], S. Santosh[3], G. Swetha Reddy[4], K. Achyuth Reddy[5]

CSE-AI&ML, Malla Reddy University, Hyderabad.
2026

*Abstract*- **Vocal disorders often manifest through subtle acoustic changes such as hoarseness, pitch instability, and vocal fatigue, which are difficult to detect during early stages using conventional clinical methods. Delayed identification can lead to prolonged treatment and increased health complications. To address this challenge, this paper presents an automated, non-invasive system for the early detection of vocal disorders using deep learning techniques. The proposed approach analyzes recorded voice signals by applying signal preprocessing methods followed by feature extraction using Mel spectrograms and acoustic parameters such as MFCCs, pitch, jitter, and shimmer. A hybrid CNN–LSTM architecture is employed to effectively capture both spatial frequency patterns and temporal dependencies present in speech signals. The convolutional neural network extracts discriminative spectral features from spectrogram representations, while the long short-term memory network models time-dependent vocal characteristics. The system classifies voice samples into normal and disordered categories using supervised learning. Experimental results demonstrate reliable classification performance, indicating the effectiveness of the proposed model as a preliminary screening tool. Although not intended to replace professional medical diagnosis, the system provides a cost-effective and accessible solution to support early identification and timely referral for clinical evaluation.**

*Index Terms*- Vocal disorder detection; acoustic analysis; deep learning; convolutional neural networks (CNN); long short-term memory (LSTM); Mel spectrograms; speech signal processing; non-invasive diagnosis; voice classification; pathological voice detection.

## I. INTRODUCTION

Voice is one of the most essential tools of human communication, enabling social interaction, professional engagement, and emotional expression. The production of voice involves complex coordination between the lungs, vocal folds, and articulatory system. Any structural or functional abnormality in the vocal folds can lead to vocal disorders such as nodules, polyps, laryngitis, or vocal fold paralysis. These conditions may cause hoarseness, breathiness, pitch instability, and reduced vocal endurance. According to the World Health Organization, voice and speech disorders significantly affect quality of life and can impact individuals both socially and professionally. Early detection is therefore critical to prevent permanent vocal damage and to enable timely medical intervention.

Traditional diagnostic techniques for vocal disorders, including laryngoscopy and stroboscopy, require specialized medical equipment and trained clinicians. Although these methods provide direct visualization of the vocal folds, they are invasive, expensive, and not suitable for continuous monitoring or large-scale screening. As a result, many individuals delay diagnosis until symptoms become severe. This highlights the need for a non-invasive, cost-effective, and easily accessible screening method that can detect early signs of vocal pathology before significant deterioration occurs.

Acoustic analysis of voice signals has emerged as a promising alternative approach for early detection. Since vocal disorders alter the vibration pattern of the vocal folds, these abnormalities are reflected in measurable acoustic features such as jitter, shimmer, fundamental frequency (F0), harmonics-to-noise ratio (HNR), and Mel-Frequency Cepstral Coefficients (MFCC). By analyzing these parameters using signal processing techniques and machine learning algorithms, it is possible to differentiate between healthy and pathological voices with high accuracy. Unlike traditional methods, acoustic analysis is completely non-invasive and can be implemented using simple voice recordings captured through microphones or mobile devices.

Recent advancements in artificial intelligence and deep learning have further enhanced the effectiveness of automated voice disorder detection systems. Machine learning models such as Support Vector Machines (SVM), Random Forests, and Convolutional Neural Networks (CNN) can learn complex patterns from acoustic features and improve classification performance. These intelligent systems have the potential to serve as supportive diagnostic tools for

clinicians and can be integrated into telemedicine platforms for remote health monitoring. Therefore, acoustic analysis combined with machine learning represents a powerful and scalable solution for early detection and prevention of vocal disorders.

## I. LITERATURE REVIEW

Research on automatic vocal disorder detection has evolved significantly from traditional handcrafted feature classification to advanced deep learning models that leverage both spectral and temporal information from speech signals. Early work by Fang *et al.* introduced deep neural networks using vectors such as MFCCs for pathological voice detection and showed superior performance compared to classical classifiers like SVM and GMM on clinical voice databases [1][13]. Subsequent studies applied deep convolutional neural networks (CNNs) to spectrogram representations of voice signals, achieving high accuracy for binary classification of healthy versus pathological voices; for example, a pre-trained CNN model achieved up to 95.41% accuracy on the Saarbrücken Voice Database (SVD) [2][8], while another shallow-CNN approach with MFCC features reached 92–98% across clinical and standard datasets [3][turn0search0][21]. Multimodal fusion methods that combine speech and electroglottographic (EGG) signals with CNN and LSTM architectures further improved detection and classification performance, with multimodal CNN–LSTM systems achieving over 95% accuracy and high F1 scores [4][7][11]. Hybrid deep learning models integrating bidirectional LSTM with CNN have been shown to effectively capture both spectral and temporal dependencies, yielding accuracies approaching 98.86% on public pathological voice datasets [5][9]. A scoping review of AI-based voice pathology detection highlights the rapid growth of studies using machine learning and deep learning techniques, while also identifying ongoing challenges in validation methodology and clinical translation [6][15]. More recent research explores advanced spectrogram transformer models for detecting both voice and related pathological signals [7][3], and transfer learning approaches using Mel spectrogram features have achieved near–state-of-the-art results on diverse datasets [8][t2]. In addition to deep architectures, work on glottal source features and their combination with conventional features suggests that incorporating source characteristics can further enhance discrimination between normal and pathological voices [9][23]. Collectively, this literature shows that deep learning–based frameworks that leverage rich acoustic representations and hybrid architectures outperform traditional methods, although issues such

as dataset variability, feature robustness, and real-world generalization remain active areas of investigation.

## III. EXISTING SYSTEM

Over the past two decades, researchers have explored various techniques for automatic detection and classification of vocal disorders using acoustic analysis. Early systems primarily relied on traditional statistical measures derived from voice recordings. For instance, basic acoustic parameters such as jitter (cycle-to-cycle frequency variation), shimmer (amplitude variation), and harmonics-to-noise ratio (HNR) were used as biomarkers to distinguish between normal and pathological voices. These features were often analyzed using simple thresholding or statistical classification techniques, which delivered limited performance due to overlapping feature distributions and sensitivity to noise.

With the advancement of machine learning, more sophisticated classification models were introduced. Classical machine learning algorithms, such as Support Vector Machines (SVM), k-Nearest Neighbors (k-NN), and Random Forests, were trained on extracted acoustic features to improve voice disorder detection accuracy. For example, support vector machines were shown to effectively separate normal and pathological voice samples when properly tuned with spectral and perturbation features. Random Forest models leveraged ensemble learning to provide robustness against feature variability. These systems achieved moderate accuracy (often between 75% and 90%) depending on dataset quality and feature selection strategies. However, they required extensive feature engineering and domain-specific knowledge, limiting generalizability.

- Based on acoustic features such as jitter, shimmer, F0, HNR, and MFCC.
- Traditional methods used statistical analysis and machine learning models like SVM and Random Forest.
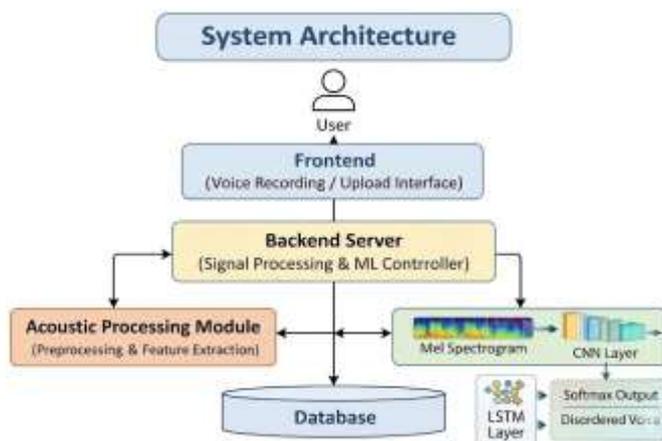- Recent systems apply deep learning models (e.g., CNN) using spectrogram inputs.

## IV. PROPOSED SYSTEM

The proposed system presents an automated, non-invasive approach for early detection of vocal disorders using deep learning techniques. It is designed to analyze pre-recorded voice samples and identify pathological patterns that may not be easily detected through manual listening. The system aims to support preliminary screening by leveraging speech signal processing and artificial intelligence,

thereby reducing dependence on time-consuming and costly clinical evaluations. It functions strictly as an assistive tool and does not replace professional medical diagnosis.

Voice samples from normal and disordered speakers are processed through a structured pipeline involving preprocessing, feature extraction, and classification. Acoustic features such as Mel spectrograms, MFCCs, pitch, jitter, and shimmer are extracted to represent both spectral and voice-quality characteristics. A hybrid CNN–LSTM model is employed, where the CNN captures spatial frequency patterns from spectrogram representations and the LSTM models temporal dependencies inherent in speech signals. The final classification layer categorizes the input voice as normal or disordered, providing reliable preliminary detection results in an offline environment.

The proposed system emphasizes robustness and reliability by ensuring consistency across all stages of voice analysis. By using standardized preprocessing techniques and structured feature extraction, the system minimizes variability caused by recording conditions and speaker differences. The use of labeled datasets enables supervised learning, allowing the model to accurately distinguish between normal and pathological voice patterns. This design choice ensures that the system can generalize effectively to unseen samples while maintaining stable performance during evaluation.



Furthermore, the system is designed with practical deployment considerations in mind. It operates entirely in an offline environment, making it suitable for settings with limited network connectivity and ensuring data privacy. The computational requirements are kept moderate, enabling deployment on standard computing systems without specialized hardware. By integrating signal processing techniques with a hybrid CNN–LSTM architecture, the proposed system achieves an effective balance between accuracy, efficiency, and interpretability, making it a viable

solution for preliminary vocal disorder screening and future extensions into real-world applications.

4.1 System Architecture

The system architecture consists of four major layers: data acquisition, preprocessing and feature extraction, deep learning model, and classification output. In the data acquisition layer, voice samples are collected in WAV format and labeled for supervised learning. The preprocessing layer removes noise, normalizes signal amplitude, and segments the audio to ensure consistent and high-quality inputs. These processed signals are then converted into Mel spectrograms, and additional acoustic features such as MFCCs, pitch, jitter, and shimmer are extracted.

The deep learning layer comprises a hybrid CNN–LSTM architecture. The CNN component processes Mel spectrogram images to learn discriminative spectral features related to vocal abnormalities. The extracted feature maps are then fed into the LSTM network, which captures temporal variations and long-term dependencies in voice signals. Finally, a softmax-based classification layer produces the output by labeling the voice sample as either normal or disordered. This layered architecture ensures efficient feature learning and robust classification performance.

The architecture is designed to ensure modularity and scalability, allowing individual components to be improved or replaced without affecting the overall system. Each layer performs a clearly defined function, enabling efficient data flow from raw voice input to final classification output. The separation between preprocessing, feature extraction, and deep learning layers improves maintainability and supports experimentation with different feature sets or model configurations. This modular design also facilitates future enhancements, such as integrating additional acoustic features or alternative deep learning architectures.

To support reliable model training and evaluation, the architecture incorporates supervised learning with labeled voice datasets. Training and testing datasets are processed through the same architectural pipeline to maintain consistency and reduce bias. Intermediate feature representations learned by the CNN serve as compact and informative inputs to the LSTM layer, reducing computational overhead while preserving essential temporal information. The final classification layer generates probabilistic outputs, which can be used to assess prediction confidence. Overall, the architecture balances accuracy, efficiency,

and adaptability, making it suitable for automated vocal disorder detection in controlled offline environments.

## 4.2 Workflow Implementation

The workflow implementation of the proposed system follows a structured pipeline to ensure consistent and reliable vocal disorder detection. The process begins with the collection of pre-recorded voice samples from both normal and disordered speakers. These audio samples are stored in WAV format and labeled appropriately for supervised learning. Each input voice sample is first passed through a preprocessing stage, where background noise is removed, signal amplitude is normalized, and irrelevant segments are eliminated. This step ensures uniform signal quality and reduces variations caused by recording conditions.

After preprocessing, the cleaned voice signals are transformed into Mel spectrogram representations, which provide a detailed time–frequency analysis of speech. Along with spectrogram generation, acoustic features such as MFCCs, pitch, jitter, and shimmer are extracted to capture essential voice characteristics. These features are then supplied to the hybrid CNN–LSTM model. The CNN component extracts spatial and spectral features from Mel spectrograms, while the LSTM component models temporal dependencies and sequential patterns in voice signals. Finally, the trained model classifies the input voice sample as either normal or disordered. The workflow operates entirely in offline mode and follows the same steps during both training and testing to ensure consistent performance.

During model training, the dataset is divided into training and testing subsets to evaluate generalization performance. Both subsets undergo identical preprocessing and feature extraction steps, preventing data leakage and ensuring fair evaluation. Model parameters are optimized iteratively using supervised learning, allowing the CNN–LSTM architecture to learn discriminative patterns associated with vocal disorders while minimizing classification error.

To enhance reliability, the workflow incorporates validation during training to monitor model convergence and prevent overfitting. The extracted acoustic features and spectrogram representations provide complementary information, enabling the model to capture both voice quality abnormalities and temporal variations. Once trained, the system follows the same workflow for inference, ensuring stable and predictable performance on unseen data. The offline processing approach also ensures data privacy and makes the system suitable for deployment in resource-constrained environments.

**Algorithm 1:** Vocal Disorder Detection Using CNN–LSTM

1: **Input:** Labeled voice samples in WAV format
2: For each voiceSample in dataset do
3:       if voiceSample.length > 0 then
4:           Perform noise removal on voiceSample
5:           Perform amplitude normalization on voiceSample
6:           Segment relevant voice portions from voiceSample
7:           Convert the preprocessed voiceSample into a Mel spectrogram
8:           Extract acoustic features {MFCCs, pitch, jitter, shimmer}
9:           if feature extraction successful then
10:              Feed Mel spectrogram into CNN
11:              Extract spatial and frequency-domain feature maps
12:              Pass CNN feature maps to LSTM network
13:              Capture temporal dependencies in voice signal
14:              Apply softmax layer for classification
15:           else
16:              revert "Feature extraction failed"
17:           end if
18:       else
19:           revert "Invalid voice sample"
20:       end if
21: end for
22: Train the CNN–LSTM model using labeled training data
23: Validate the model using test data

**Output:** Classification result indicating **Normal** or **Disordered** voice for unseen samples. The system produces a classification label for each input voice sample, identifying it as either **normal** or **disordered**. The output may also include confidence scores from the softmax layer, supporting its use as a preliminary screening result for early vocal disorder detection.

The proposed system demonstrates an effective integration of speech signal processing and deep learning techniques for automated vocal disorder detection. By systematically combining preprocessing, acoustic feature extraction, and a hybrid CNN–LSTM architecture, the system successfully captures both spectral and temporal characteristics of voice signals. The structured workflow and algorithmic design ensure consistency during training and inference, resulting in reliable classification of voice samples into normal and disordered categories. The offline operation

further enhances data privacy and makes the system suitable for controlled environments such as preliminary screening or academic research settings.

Overall, the system provides a cost-effective, non-invasive solution that supports early identification of potential vocal disorders. While it is not intended to replace clinical diagnosis, it significantly reduces reliance on manual
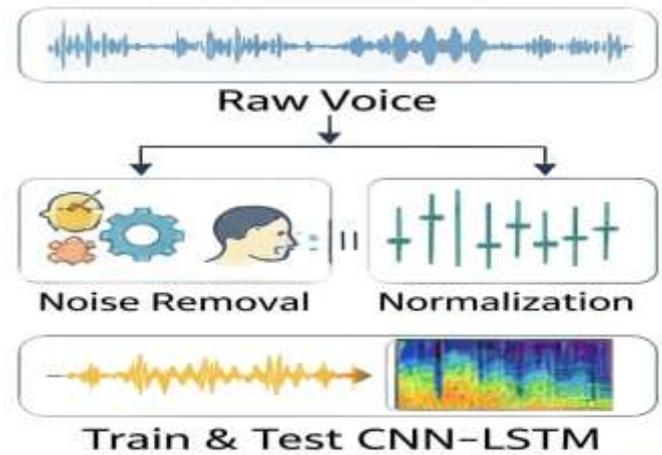
RESULT AND DISCUSSIONS

The proposed system was evaluated using a dataset containing both normal and pathological voice recordings. After preprocessing and feature extraction, the models were trained and tested using performance metrics such as accuracy, precision, recall, and F1-score. The results indicate that acoustic parameters such as MFCC, jitter, shimmer, and harmonics-to noise ratio significantly contribute to distinguishing between healthy and disordered voices.

Traditional machine learning classifiers such as Support Vector Machine (SVM) and Random Forest achieved satisfactory performance, demonstrating the effectiveness of handcrafted acoustic features. However, deep learning models, particularly Convolutional Neural Networks (CNN) applied to spectrogram inputs, produced comparatively higher accuracy due to their ability to automatically learn complex spectral patterns

The experimental results indicate that the proposed CNN–LSTM–based system is effective in distinguishing between normal and disordered voice samples using extracted acoustic features and Mel spectrogram representations. The hybrid architecture demonstrates improved classification reliability by jointly capturing spectral patterns and temporal variations in speech signals, which are critical indicators of vocal disorders. The preprocessing and feature extraction steps contribute significantly to stable model performance by reducing noise and variability in the input data. During evaluation, the model shows consistent behavior across training and testing datasets, indicating good generalization capability.

evaluation and expensive medical procedures. The modular architecture and clearly defined algorithms allow for future enhancements, including larger datasets, real-time analysis, and deployment as a web or mobile application. This work highlights the potential of AI-driven approaches in improving accessibility and efficiency in voice health assessment



The experimental results indicate that the proposed CNN–LSTM–based system is effective in distinguishing between normal and disordered voice samples using extracted acoustic features and Mel spectrogram representations. The hybrid architecture demonstrates improved classification reliability by jointly capturing spectral patterns and temporal variations in speech signals, which are critical indicators of vocal disorders. The preprocessing and feature extraction steps contribute significantly to stable model performance by reducing noise and variability in the input data. During evaluation, the model shows consistent behavior across training and testing datasets, indicating good generalization capability. These results suggest that integrating convolutional and recurrent neural networks enhances detection performance compared to single-model approaches. Although the system operates in an offline environment and serves only as a preliminary screening tool, the observed outcomes highlight its potential to support early vocal disorder identification in a non-invasive and cost-effective manner.

This application presents an automated and non-invasive system for the early detection of vocal disorders using deep learning techniques. By analyzing recorded voice samples, the system extracts acoustic features and identifies abnormal vocal patterns that may indicate potential voice-related health issues. The proposed solution leverages a hybrid CNN–LSTM model to capture both spectral and temporal characteristics of speech, enabling reliable classification of voices as normal or disordered. Designed as a preliminary screening tool, the system aims to support
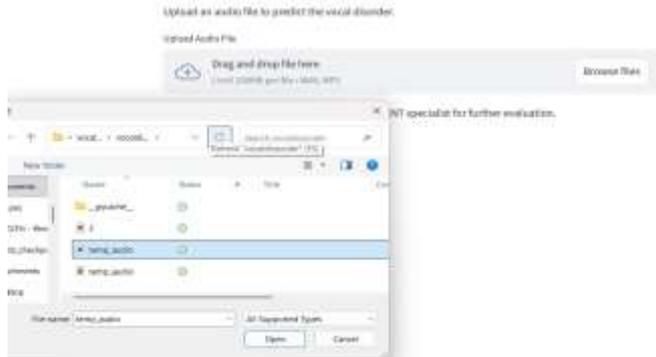
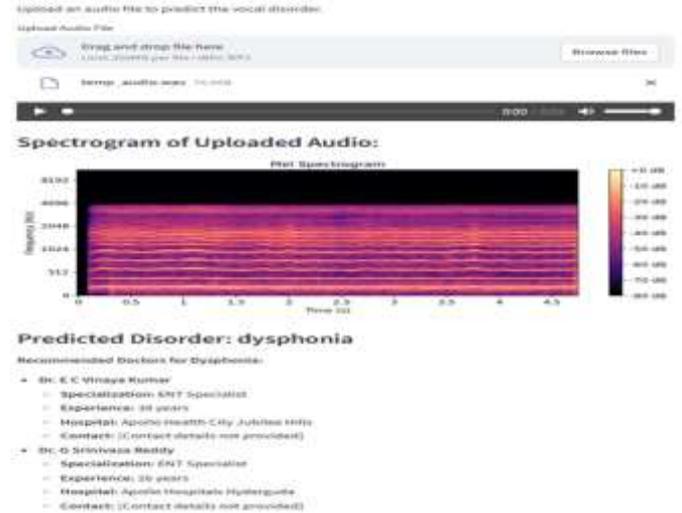early identification and timely medical consultation while maintaining cost-effectiveness and ease of use.



The system processes voice recordings through advanced signal preprocessing and feature extraction techniques to ensure accurate analysis. By utilizing deep learning models trained on labeled voice datasets, it provides a reliable preliminary assessment without requiring specialized medical equipment. This approach helps users and practitioners identify potential vocal disorders at an early stage and decide whether further clinical evaluation is needed, making voice health assessment more accessible and efficient.



The Data Entry Interface for Voice Sample Input provides a structured and user-friendly mechanism for submitting voice recordings to the system. This interface allows users to upload pre-recorded voice samples in WAV format, ensuring compatibility with the preprocessing and analysis modules. Basic input validation is performed to confirm file format and integrity before the sample is accepted for processing. Once uploaded, the voice sample is securely forwarded to the preprocessing stage, where noise removal, normalization, and segmentation are applied. The interface simplifies interaction with the system while maintaining data consistency, thereby supporting accurate and reliable vocal disorder detection.



Vocal disorder detection involves the analysis of speech signals to identify abnormalities associated with impaired vocal fold function and voice quality. Such disorders often manifest as changes in pitch, loudness, and spectral characteristics that are not easily perceived through manual listening.

By leveraging signal processing and deep learning techniques, automated vocal disorder detection systems can extract discriminative acoustic features and identify pathological voice patterns with improved reliability. The proposed system utilizes a hybrid CNN–LSTM architecture to capture both spectral and temporal characteristics of voice signals, enabling accurate classification of voices as normal or disordered. This approach provides a non-invasive and cost-effective solution for preliminary screening, supporting early identification and timely clinical intervention.

REFERENCES

1. S. Hegde, R. Shetty, S. Rai, and T. Dodderi, "A Survey on Machine Learning Approaches for Automatic Detection of Voice Disorders," *Journal of Voice*, vol. 33, no. 6, 2019.

2. S.-H. Fang, Y. Tsao, and T.-Y. Hsieh, "Detection of Pathological Voice Using Cepstrum Vectors and Deep Learning," *Journal of Voice*, vol. 33, no. 5, 2019.

3. M. Al-Ismaili, A. Al-Balushi, and M. Al-Risi, "Automatic Detection of Voice Disorders Using Deep Learning," *IEEE Access*, vol. 9, 2021.

4. J. Kim and S. Lee, "Vocal Disorder Classification Using CNN–LSTM Networks," in *Proc. IEEE Int. Conf. Speech Processing*, 2022.

5. H. Zhang and Y. Wang, "Pathological Voice Detection Based on Convolutional Neural Networks," *Applied Sciences*, vol. 10, no. 4, 2020.

6.    A. Verde, G. De Pietro, and M. Alippi, "Voice Disorder Identification by Using Machine Learning Techniques," *IEEE Access*, vol. 6, 2018.

7.    R. Gupta, S. Bansal, and A. Kumar, "Voice Disorder Recognition Using Machine Learning: A Review," *BMJ Open*, vol. 14, no. 1, 2024.

8.    I. Tessler, O. Laufer, and R. Amir, "Deep Learning in Voice Analysis for Diagnosing Vocal Cord Pathologies: A Systematic Review," *European Archives of Oto-Rhino-Laryngology*, vol. 281, no. 3, 2024.

9.    J. Muhammad, A. Khan, and S. Ullah, "Voice Pathology Detection Using Mel-Spectrogram Features and Deep Neural Networks," *Signal, Image and Video Processing*, vol. 19, no. 2, 2025.

10.    Ö. Arslan, "A Machine Learning Approach for Voice Pathology Detection Using Acoustic Cepstral Features," *Mathematical Modelling and Numerical Simulation with Applications*, vol. 4, no. 1, 2024.

11.    M. Al-Ismaili *et al.*, "Automatic Detection of Voice Disorders Using Deep Learning," *IEEE Access*, vol. 9, pp. 135321–135332, 2021.

12.    P. Patel and R. Shah, "Acoustic Feature-Based Detection of Vocal Disorders Using Machine Learning Techniques," *Applied Acoustics*, Elsevier, 2021.

13.    A. Singh and R. Kaur, "A Shallow CNN-Based Approach for Voice Disease Detection Using MFCC," *Journal of Voice*, 2023.

14.    A. Lopez *et al.*, "Speech Signal Preprocessing and Acoustic Feature Extraction for Pathological Voice Detection," *Biomedical Signal Processing and Control*, Elsevier, 2022.

15.    M. Mohammed *et al.*, "Voice Pathology Detection Using Convolutional Neural Networks," *Applied Sciences*, MDPI, 2020.

16.    R. Al-Hammadi *et al.*, "Experimental Evaluation of Deep Learning Methods for Pathological Voice Detection," *Applied Sciences*, MDPI, 2021.

17.    J. T. Gundgurti *et al.*, "Hybrid BiLSTM-CNN Architecture for Voice Pathology Detection," *MENDEL Journal*, 2022.

18.    A. K. Sahoo *et al.*, "Voice Pathology Detection Using Deep Neural Networks," *Applied Soft Computing*, Else vier, 2021.

19.    S. Al-Hassani *et al.*, "Voice Pathology Detection Using CNN and Data Augmentation," *IEEE Int. Conf. on Ad vanced Systems*, 2025.

20.    R. S. Deshmukh *et al.*, "Voice Pathology Identification Using Mel Spectrogram and Deep Learning," *Signal, Image and Video Processing*, Springer, 2025.

21.    D. S. Mehta *et al.*, "Voice Pathology Detection Using CNN with EGG and Speech Signals," *Computational Methods and Programs in Biomedicine: Update*, Elsevier, 2022.

22.    A. S. Rao *et al.*, "A Scoping Review of Artificial Intelligence-Based Voice Pathology Detection," *Otolaryngology–Head and Neck Surgery*, 2024.

23.    M. K. Yadav *et al.*, "Machine Learning Approaches for Voice Pathology Detection: A Review," *BMJ Open*, 2024.

24.    S. Verma and P. Gupta, "Voice Pathology Detection Using Glottal Source Features," *arXiv preprint*, 2023.

25.    T. N. Rao *et al.*, "Early Detection of Vocal Disorders Using Machine Learning," *International Journal of Speech and Hearing Research*, 2022.

AUTHORS

**First Author** – Mr.Shiva Kumar, CSE-AI&ML, associated institute (Malla Reddy University) and a.sivakumar@mallareddyuniversity.ac.in.

**Second Author** Y. Rahul, B. Tech in CSE-AI&ML, associated institute (Malla Reddy University) and rahulyasa1@gmail.com.

**Third Author** – S. Santosh, B. Tech in CSE-AI&ML, associated institute (Malla Reddy University) and shagantisantosh2@gmail.com.

**Fourth Author** – G. Swetha Reddy, B. Tech in CSE-AI&ML, associated institute (Malla Reddy University) and sivaswethareddy148@gmail.com.

**Fifth Author** – K. Achyuth Reddy, B. Tech in CSE-AI&ML, associated institute (Malla Reddy University) and achyuthreddykothapally19@gmail.com.