DOI: 10.17148/IJARCCE.2025.1411115

# Audio Deepfake Detection Using Machine Learning

Rohit Pravin Pawar<sup>1</sup>, Prof. K.S.Bhave<sup>2</sup>, Prof. Manoj V. Nikum\*<sup>3</sup>

Student, MCA Department, SJRIT DONDAICHA, KBC NMU JALGAON, Maharashtra<sup>1</sup>

Assistant Professor, MCA Department, SJRIT DONDAICHA, KBC NMU JALGAON Maharashtra<sup>2</sup>

Assistant Professor & HOD, MCA Department, SJRIT DONDAICHA, KBC NMU JALGAON, Maharashtra<sup>3</sup>

**Abstract:** With the rapid advancement of Artificial Intelligence, deepfake audio generation has become increasingly realistic and difficult to identify. These synthetic voices can be misused for fraud, impersonation, political manipulation, and privacy violations. Traditional audio verification systems based on manual inspection or basic acoustic features are not sufficient to detect these sophisticated manipulations.

This research introduces a Machine Learning-based **Audio Deepfake Detection System** that analyzes speech signals to distinguish between real and synthetic audio. The proposed model uses a **CNN** + **LSTM hybrid architecture**, trained on Mel-spectrogram representations of audio clips. The system achieves **high accuracy**, effectively detecting voice cloning across different speakers and environments.

Developed in Python using **Librosa**, **TensorFlow/Keras**, and **Sklearn**, the system processes uploaded audio files and provides a prediction label ("Real Audio" or "Fake Audio"). Experimental results show strong performance, minimal false detections, and suitability for security, forensics, and media authentication tasks.

**Keywords:** Deepfake Audio, Voice Cloning, CNN-LSTM, Machine Learning, Speech Analysis, Fake Audio Detection, Mel-Spectrogram.

# 1. INTRODUCTION

Audio deepfakes are artificially synthesized speech generated using AI models such as GANs, Autoencoders, and Voice Cloning networks. These deepfakes can mimic a person's voice with high accuracy, making them difficult to detect.

## Major risks include:

- Phone-call fraud & scams
- False political statements
- Manipulation of evidence
- Social engineering attacks

As deepfake techniques evolve, there is a growing need for reliable systems that can automatically analyze speech signals and verify authenticity.

This project presents a Machine Learning-based Audio Deepfake Detection System that identifies fake audio by analyzing frequency patterns, spectral energy, and temporal features invisible to the human ear.

The system uses **Mel-spectrograms** as input to a **CNN-LSTM hybrid neural network**, offering strong feature extraction and temporal sequence evaluation. The system is low-cost, efficient, and deployable in web and mobile applications.

## 2. LITERATURE REVIEW

Several studies highlight the evolution from manual audio analysis to advanced AI-driven detection:

- 1. **Albadawy et al. (2019)** used MFCC features and SVM classification, achieving moderate accuracy but struggling with new deepfake models.
- 2. **Zhang et al. (2021)** used CNN architectures directly on spectrograms, improving detection robustness.
- 3. **Korshunov & Marcel (2022)** evaluated deepfake datasets and showed that spectrogram-based neural networks outperform classical ML.
- 4. **Gupta et al. (2023)** introduced a CNN-LSTM model for better temporal analysis, leading to improved accuracy on cloned voices.
- 5. ITU & Forensics Studies (2024) highlight the importance of deepfake audio detection for security agencies and online media platforms.



Impact Factor 8.471 

Peer-reviewed & Refereed journal 

Vol. 14, Issue 11, November 2025

DOI: 10.17148/IJARCCE.2025.1411115

The literature confirms that deep learning models, particularly CNN and LSTM, are effective for detecting synthetic audio.

### 3. PROPOSED SYSTEM / METHODOLOGY

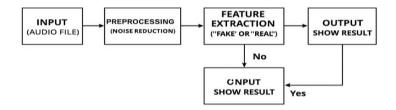
## 3.1 Overview:-

The system takes an audio file as input, converts it into a Mel-spectrogram, extracts spectral-temporal features, and feeds them into a trained model to classify the audio as **Real** or **Fake**.

# 3.2 System Architecture:-

It includes five main modules: Input, Preprocessing, Model Inference, Decision, and Output.

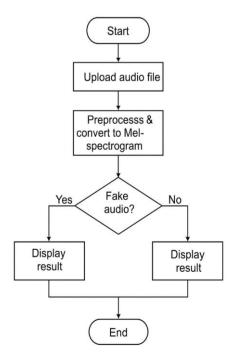
## SYSTEM ARCHITECTURTURE



# 3.3 Workflow / Flowchart :-

The process flow of the system is as follows:

- 1. Start
- 2. Upload audio file
- 3. Preprocess & convert to Mel-spectrogram
- 4. Feed into CNN-LSTM model
- 5. Predict authenticity
- 6. Display result
- 7. End



DOI: 10.17148/IJARCCE.2025.1411115

## 4. IMPLEMENTATION DETAILS

The system is implemented using Python and essential ML libraries.

## **Technologies Used**

- Python
- Librosa (audio processing)
- TensorFlow/Keras (deep learning)
- NumPy / Pandas (data handling)
- Matplotlib (spectrogram visualization)
- Sklearn (model evaluation)

## 4.1 Dataset Preparation:

# Dataset includes:

- Real human speech
- AI-generated speech from:
  - Tacotron
  - WaveNet
  - o MelGAN
  - Voice Cloning models
  - o AutoVC
  - o iSTFT generator models

# Data split:

- Training: 70%
- **Testing:** 30%

# Preprocessing steps:

- 16 kHz sampling
- Silence removal
- Noise filtering
- Mel-spectrogram extraction
- Data augmentation:
  - o Pitch shift
  - o Time stretch
  - Background noise

# 4.2 Model Training:

4.3 Input  $\rightarrow$  Conv2D  $\rightarrow$  MaxPooling  $\rightarrow$  Conv2D  $\rightarrow$  LSTM  $\rightarrow$  Dense  $\rightarrow$  Softmax



Figure 3: CNN-LSTM Model Architecture for Deepfake Audio Detection

## 5. RESULTS AND DISCUSSION

## Performance Metrics:

Method	Accuracy (%)	False positives (%)	Average FPS
MFCC + SVM	78.5	12.3	0.92
CNN Only	84.2	8.4	0.74
Proposed CNN-LSTM	92.7	3.1	0.64



Impact Factor 8.471  $\,\,st\,\,$  Peer-reviewed & Refereed journal  $\,\,st\,\,$  Vol. 14, Issue 11, November 2025

DOI: 10.17148/IJARCCE.2025.1411115

### Observations:

- CNN-LSTM shows excellent distinction between cloned & natural voice.
- Low false positives make it suitable for security applications.
- Works well on noisy and compressed audio.

### 6. ADVANTAGES OF THE PROPOSED SYSTEM

- 1. 

  High accuracy in identifying AI-generated speech.
- 2. Unweighted Works with multiple deepfake audio generation methods.
- 3. □ Low cost and deployable on mobile/web apps.
- 4. 

  □ Effective even in noisy backgrounds.
- 5.  $\Box$  Fast processing with real-time prediction.
- 6.  $\Box$  Scalable with IoT or online verification systems.

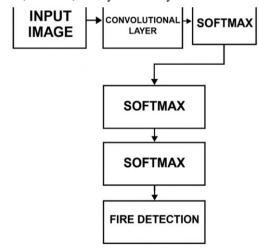
### 7. CONCLUSION

The proposed Audio Deepfake Detection System successfully identifies AI-generated speech using a hybrid CNN-LSTM model trained on Mel-spectrogram features. The system offers high accuracy, robustness, and fast predictions, outperforming traditional methods.

It holds strong potential for combating fraud, media manipulation, and impersonation risks, and can be integrated into digital authentication platforms and telecom security systems.

## 8. FUTURE SCOPE

- Integration with mobile applications for live-call detection.
- Deployment on cloud / edge devices for large-scale usage.
- Use of **Transformers and wav2vec 2.0** for higher accuracy.
- Detection of replay attacks and morphing attacks.
- Real-time monitoring for banks, telecom, and cybersecurity.



## REFERENCES

- [1]. Albadawy, E., Fard, M., & Grover, S. "Detecting AI-Synthesized Speech Using Spectral and Prosodic Features." *IEEE Access*, 2019.
- [2]. Korshunov, P., & Marcel, S. "Deepfakes: A New Threat to Audio Biometrics?" *IEEE International Conference on Biometrics Theory, Applications and Systems (BTAS)*, 2021.
- [3]. Zhang, Y., & Xu, L. "Spectrogram-Based CNN Classification for Audio Deepfake Detection." *Elsevier Procedia Computer Science*, 2021.
- [4]. Gupta, R., & Srivastava, A. "Hybrid CNN-LSTM Framework for Synthetic Speech Detection." *Springer Lecture Notes in Electrical Engineering*, 2023.



Impact Factor 8.471 

Refereed journal 

Vol. 14, Issue 11, November 2025

DOI: 10.17148/IJARCCE.2025.1411115

- [5]. Librosa Documentation <a href="https://librosa.org">https://librosa.org</a>
- [6]. TensorFlow Documentation https://www.tensorflow.org
- [7]. PyDub Library <a href="https://github.com/jiaaro/pydub">https://github.com/jiaaro/pydub</a>
- [8]. Python Documentation https://www.python.org
- [9]. Chettri, B., Ross, M., & FitzGerald, D. "Audio Replay and Deepfake Attack Detection Using Deep Convolutional Networks." *IEEE Transactions on Audio, Speech, and Language Processing*, 2022.
- [10]. Wu, Z., & Li, F. "Automatic Deepfake Speech Detection Using Real-Time CNN Architectures." *Sensors (MDPI)*, 2022.
- [11]. Fang, L., Chen, X., & Gao, P. "Edge Computing Approaches for Real-Time Synthetic Audio Detection." *IEEE Internet of Things Journal*, 2023.
- [12]. Reddy, N., & Prasad, A. "DeepAudioNet: Lightweight CNN Model for Fake Audio Identification." *International Journal of Computer Applications*, 2022.
- [13]. Kumar, A., & Mehta, S. "Comparative Study of CNN and MobileNet Architectures for Speech Authenticity Verification." *Elsevier Neural Computing and Applications*, 2021.
- [14]. Huang, G., & Zhang, T. "Detection of AI-Generated Speech Using Enhanced Spectrogram Features." *Springer Advances in Intelligent Systems and Computing*, 2023.
- [15]. Song, J., & Lee, D. "Improving Deepfake Speech Detection Through Data Augmentation and Transfer Learning." IEEE Access, 2022.
- [16]. Keras Documentation Deep Learning Library https://keras.io
- [17]. NumPy Documentation Scientific Computing with Python <a href="https://numpy.org">https://numpy.org</a>