



# IJRASET

International Journal For Research in  
Applied Science and Engineering Technology



# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

**Volume:** 12    **Issue:** XI    **Month of publication:** November 2024

**DOI:** <https://doi.org/10.22214/ijraset.2024.64682>

[www.ijraset.com](http://www.ijraset.com)

Call:  08813907089

E-mail ID: [ijraset@gmail.com](mailto:ijraset@gmail.com)

# Speaker Verification Model USING LST

Prof. Supriya Telsang<sup>1</sup>, Pratik Pawar<sup>2</sup>, Prachi Pawar<sup>3</sup>, Krushna Pawar<sup>4</sup>, Harish Pawar<sup>5</sup>, Darshan Pawar<sup>6</sup>  
Department of Engineering, Sciences and Humanities (DESH) Vishwakarma Institute of Technology, Pune, India

**Abstract:** *This project presents a speaker verification model using LSTM networks to detect the authenticity of voices, distinguishing between real and deepfake audio. With the growing threat of AI-generated voices being used for identity theft and misinformation, this model aims to address those risks. The system uses MFCCs (Mel-frequency cepstral coefficients) to capture key speech features and distinguish real voices from altered ones. The dataset includes both real and deepfake samples from well-known public figures like Joe Biden and Donald Trump, with techniques like time-stretching and pitch-shifting used to improve the model's performance.*

*Built with LSTM layers and trained for 50 epochs using the Adam optimizer, the model achieved an accuracy of 84.62%, proving effective in detecting fake voices. This LSTM-based model offers a valuable solution for enhancing voice authentication systems, especially in critical areas like security systems, online transactions, and biometric authentication. Its robust performance demonstrates the potential to counter the rising threat of deepfake technology.*

**Keywords:** *LSTM networks, MFCC, Voice Cloning Detection Real-Time Speech Verification, Adam Optimizer for Audio Models, Machine Learning for speech Security, Machine learning, security system*

## I. INTRODUCTION

Over the years, the capabilities of artificial intelligence develop to an extent where technological advancements allow the duplication of voices and real time transformation of identities. It is an impressive technological advancement, but the use of such sophisticated technology comes with ethical and safety risks. Being able to make deepfake voices is an opportunity that can be exploited through identity theft, scams, and impersonation, especially in areas where voice authentication is used, and in general communications. Speaker verification is one of the techniques, commonly used in modern banking systems, security systems, and voice interactive devices, to confirm an individual's identity using voice. Nevertheless, such systems are in danger of being enhanced by deepfake voices that readily reproduce anyone's voice and hence render the traditional speaker verification models ineffective. Such models hardly differentiate between original human voice and artificially produced voice. As a solution to these problems, we develop a novel speaker verification system that is able to discriminate between real and deep-fake speech. Our model is based on Long Short Term Memory (LSTM) networks which will help understand the temporal correlations existing in speech data. Additionally, we include Mel-frequency cepstral coefficients (MFCCs) which is a technique widely incorporated in speech recognition, and it helps in capturing important aspects of the voice which aid in identifying real from fake samples. For this project, the developer used a dataset, called the DEEP-VOICE dataset, which consisted of real and synthesized audio samples of well-known personalities such as Joe Biden and Donald Trump. Moreover, we used time-stretching and pitch-shifting methods as data augmentation strategies to increase the variation of the training data and improve the performance of the model. After training the model for 50 epochs with the Adam optimizer, the model attained a reaches a test accuracy of 84.62%, which could classify the voices as real or fake. This a model is quite effective in addressing the issue of voice authentication systems and their limitations especially in sensitive areas such as security, online transaction, and voice biometrics. By proposing a solution to the threat of deepfake technology, this research work aims to improve the systems' trustworthiness and the overall security of speaker verification systems.

### A. Background

As technology, voice recognition and speaker verification have become integral components of modern security systems. These technologies allow for the identification and authentication of individuals based on their voice patterns, which are unique to each person. Speaker verification is widely used to each person. Speaker verification is widely used in various applications, including banking, telecommunications, home security, and biometric systems. This method of identification is both convenient and accessible, as it allows users to verify their identity without the need for physical credentials, such as passwords or fingerprints.

Nonetheless, new advancements in artificial intelligence (AI), especially in its use for speech production, have emerged a difficult issue: the ability to produce fictitious deepfake voices. Using deepfake technology, it becomes possible to create an artificial voice that mimics the vocal characteristics of an actual specific person. Such voices can be made with AI so advanced that it can be virtually impossible to tell whether any particular voice is real or fake for the purposes of speaker verification.

Such deepfake voices are created through many ways, one of them being retrieval-based voice conversion, which essentially means modifying a person's voice to match a certain other voice. This technology is effective in many sectors including movies and devices for specialized assistance to the disabled, but also has a negative aspect that is a serious danger to privacy and security. Criminals are also able to use such deceptive voice creations in order to sound like a person and make someone let them in a secure system.

Owing to the above delineated concerns, there is a high demand for sophisticated systems that have the ability to detect deepfake audio in a very short time. Existing speaker verification systems were mainly created to authenticate only real person's voices, and thus do not stand the advanced forgeries which the deepfake technology creates. The Challenges in Identifying Speech Synthesizer Technology in Real-Time therefore becomes an essential research direction, especially where voice recognition system is employed in putting security measures.

In this research, we seek to solve the aforementioned problems by designing a speaker verification model which is inclusive not only speech recognition but also detection of deepfake audio. Therefore, in the work presented in this paper, we attempt to develop a model which integrates LSTM networks together with cepstral features of the audio signal to enhance deepfake voice detection in real world applications.

## II. LITERATURE SURVEY

The rapid advancements in machine learning and deep learning have significantly impacted speech recognition and speaker verification technologies. In this literature survey, we explore the foundational research and state-of-the-art methods relevant to the development of a speaker verification model using Long Short-Term Memory (LSTM) networks, with a particular focus on detecting deepfake voices.

### A. Speaker Verification and Traditional Methods

Speaker verification is a biometric process that uses the unique characteristics of a person's voice to authenticate their identity. Early methods for speaker verification relied on techniques like Gaussian Mixture Models (GMMs) and Hidden Markov Models (HMMs). These statistical models worked well for basic speaker recognition tasks, but they struggled with more complex scenarios, such as background noise or varying speaking conditions.

Research by Reynolds (2000) introduced GMM-UBM (Universal Background Model), which improved the robustness of speaker verification systems. However, as voice manipulation and deepfake technologies have advanced, these traditional models have become less effective in detecting synthetic or cloned voices.

### B. Deep Learning Approaches

With the rise of deep learning, newer models such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) have gained popularity in the speaker verification domain. These models have the ability to learn complex patterns from raw audio data, enabling them to outperform traditional methods.

In recent studies, Long Short-Term Memory (LSTM) networks, a specific type of RNN, have proven to be particularly effective in capturing temporal dependencies in speech. Graves et al. (2013) demonstrated that LSTMs could model time-series data more effectively than standard RNNs, which are prone to issues like vanishing gradients. This makes LSTMs highly suitable for tasks like speech recognition and speaker verification, where the temporal context of the speech signal is critical.

### C. MFCCs in Speech Processing

Mel-frequency cepstral coefficients (MFCCs) have been widely adopted in speech and speaker recognition tasks due to their ability to capture the acoustic characteristics of speech. First introduced by Davis and Mermelstein (1980), MFCCs provide a compact and powerful representation of speech that aligns with the human auditory system.

MFCCs extract both the spectral envelope and temporal characteristics of speech, which are important for distinguishing between speakers. Numerous studies have confirmed the effectiveness of MFCCs in both traditional and deep learning-based speaker verification systems. For instance, Nagrani et al. (2017) used MFCCs in combination with deep learning models for speaker identification tasks, demonstrating strong performance in large-scale datasets like VoxCeleb.

#### D. Deepfake Detection in Speech

As deepfake technology has advanced, detecting AI-generated or cloned voices has become a critical challenge in speaker verification. Deepfake voices can be generated using Voice Conversion (VC) techniques or Text-to-Speech (TTS) systems, making it difficult for traditional speaker verification models to distinguish between real and synthetic voices.

Recent studies, such as Kinnunen et al. (2020), have explored the vulnerabilities of speaker verification systems to deepfake attacks. They propose that models designed for speech synthesis or voice recognition must evolve to incorporate deepfake detection capabilities. Researchers have also begun integrating data augmentation techniques, such as pitch-shifting and time-stretching, to improve model robustness and enhance the detection of AI-manipulated voices.

In particular, Xie et al. (2021) showed that deep learning models like LSTMs, when combined with robust features like MFCCs, can significantly improve the detection of AI-generated speech. This aligns with our approach of using LSTMs and MFCCs for building a more resilient speaker verification model capable of handling deepfake threats.

#### E. LSTM Networks in Speaker Verification

LSTM networks have been particularly effective in speaker verification tasks due to their ability to learn long-term dependencies in time-series data. LSTMs are designed to retain information across long sequences, which is crucial for understanding the structure of speech over time.

In the study by Zhang et al. (2016), LSTMs were shown to outperform traditional models in speaker verification tasks, particularly in handling variations in speech, such as different speaking rates and background noise. The gated architecture of LSTMs allows them to selectively remember important features and forget irrelevant ones, which is essential in tasks involving sequential data like speech.

By utilizing LSTM networks, our model leverages these strengths to learn the patterns of both real and synthetic voices, improving the detection of deepfake audio.

#### F. Applications in Everyday Life

Speaker verification has a wide range of applications in real-world scenarios, from securing banking transactions and personal devices to protecting sensitive communications. With the increasing prevalence of deepfake voices, it is more important than ever to build systems capable of distinguishing real voices from fake ones.

The integration of LSTM networks and MFCCs into speaker verification systems offers a promising solution to this growing threat. As deepfake detection becomes a key area of focus, future models will need to incorporate real-time verification capabilities, enabling everyday applications such as smart home devices, voice-activated assistants, and telecommunications systems to remain secure.

### III. METHODOLOGY / FINDINGS

The objective of our investigation is to create a speaker evaluation system capable of identifying whether a voice is human or artificially synthesized through the use of voice deepfake technology. In pursuing this goal, we have combined Long Short-Term Memory (LSTM) networks and Mel-frequency cepstral coefficients (MFCCs) for the audio-based processing. Below is provided an overview of the important processes in constructing the model.

#### A. Dataset Creation and Preparation

For this research, we employed the DEEP-VOICE dataset specifically developed for this consideration. It contains two categories of audio samples: authentic recordings and deliberate misrepresentation audio samples. The former ones were sourced from eight publicly recognized figures, which included political icons such as Joe Biden and Donald Trump. The latter's voice sampling utilized Retrieval-based Voice Conversion that transforms the voice of a specific person into the likeness of another individual. To ensure improvement in the model's learning on unseen data, we performed some techniques in data clear up to add diversity to the given work. We did this by using:

Time-stretching this involves elongation or shortening of audio without altering the pitch; this helps demonstrate how the speed of individuals' speech may vary.

Pitch-shifting works on a similar audio format that employs a sample which is unforgettable; however, its speed is retained, saw different voices changing tones from below to higher tones.

These techniques help the model to become more robust by helping it cope with variations in human speech which may not exactly mimic the speech used in training.

### B. Feature Extraction: MFCCs

For the analysis of the audio, we have used Mel-frequency cepstral coefficients (MFCCs) obtained from each audio sample. MFCCs are widely used in speech recognition as this parameter considers many active factors of speech, such as voice modulation, frequency, and quality of sound, which helps to differentiate one speaker from others. In the paper, we have collected 13 MFCCs from each audio sample for the purpose of analysis. These features were then used as the input data for our model.

### C. Model Architecture: LSTM Network

This model is a type Long Short-Term Memory Networks (LSTM), a specialized Recurrent Neural Network (RNN) which is appropriate where data is temporal in nature such as speech. LSTM networks are more advanced than traditional RNN networks because they can remember information for long periods, learning how various portions of the speech sample correlate with one another after some time.

The model incorporates the following components:

Two LSTMs with 64 and 32 units respectively. These layers enable the model to capture the temporal aspects of the audio.

Dropout layers with a tendency of 0.2 in between the LSTM layers. This layer reduces overfitting, that is the tendency of a model to learn specific patterns to fit the training dataset but fails to generalize to new datasets.

Dense (fully connected) layers that contribute in cubing the 'deep net' architecture and help to regard if the voice is real or fake.

A softmax layer at the end assigning probabilities to the two classes, fake or real.

### D. Training and Evaluation

Adam has been applied as the training algorithm. It is a synonymous term with optimization process during deep learning. The learning rate was established at 0.001 that is the rapidity of which how the model adjusts its parameters. A batch size of 32 was also used, that is, during training, the model trains on 32 audio readings at a time.

The duration for which the model was trained is fifty epochs which means the model has scrutinized the entire dataset fifty times over. For this purpose, we restricted the input set to 80% for training and 20% for testing the effectiveness of the algorithm. To add even more credibility to our results, we also included 10-fold cross-validation as well. This approach takes a given data set splits it into 10 segments, takes 9 of them and trains the model while testing the last 1. This is done 10 times with every segment used once as the test sample.

## IV. RESULTS AND DISCUSSIONS

The glove-based communication interpreter system has been created and examined in a variety of environments in order to test its effectiveness in a live military operation. This section provides an overview of these tests and examines the possible consequences of the system on communication on the battlefield.

### A. System Performance and Accuracy

The performance evaluation of the system was based on the accuracy of syntax and verbal commands recognition ability for selected military gestures. Test evaluation imposed various sized and dexterous users performing a series of defined gestures and the system was very accurate over 90% of the features were correctly recognized and translated. The addition of flex sensors and accelerometer allowed for more efficient and effective capturing of both dynamic hand postures and movements inputs which increases the reliability of the system.

In the training mode, the system was able to record and remember specific gestures designed for each individual user. As a result, the gesture was correctly recognized in the operational mode, even if the gesture differed slightly due to individual hand movement differences.

The system is designed in an adaptive manner and therefore can support users with a very wide spectrum making it more useful in the military where there is diversity in users.

### B. Real-Time Usage and Field Testing

The analysis of the glove performance included an assessment of its 24–7 usability under simulated battlefield situations. These tests involved the great room for motion, extreme noise, poor lighting conditions and fast pace action, which are common in war. The system performed well in these conditions, as the glove interface was able to relay commands from the commander to the troops effectively.

For instance, the hand movements of the commander were swiftly converted to text on a screen and voice commands that were then sent to the earpieces of the soldiers. This enabled soldiers to communicate without making sounds, eliminating the disturbances that could come with making oral orders. Such concerns as the excessive noise in the environment or being overheard by enemy combatants muted the use of vocal instructions.

The system enables the soldier to receive both visual and sound inputs simultaneously most importantly when instructions are given due to the urgency of the situation. The interaction of these factors enhances communication and eliminates some risks. For instance, the use of hand signals, which is also quite effective in the communication process, does not easily expose the user to the enemy.

### C. *Impact on Military Operations*

Deployment of this glove-based system will fundamentally change communication of commands in the field. Since the system supports communication through physical gestures without any spoken words, it resolves some of the issues associated with the conventional military communication practices. Additionally, since gesture recognition can be tailored, the commanders can design their conventional hand signs for missions and units which can help fortify operational security. Such customization means that the commands can only be known to the people who are supposed to interpret them, and who have mastered the gestures set, thus minimizing the chances of miscommunication or interception by enemies.

The time between an issued command and the action being taken by the soldiers is almost non-existent.

This is mostly under combat scenarios that are dynamic, where every second occurs with a decision that can cause success or failure.

### D. *Challenges and Limitations*

Even so, the system's appreciable potential does come with encumbering challenges and limitations that should be solved. Priority among these challenges is the need to keep the system functioning in extreme situations such very high and low temperatures, as well as instances where the glove would get hit. The current system operates in the limits of range and functionality of the Bluetooth module, whose performance may be hindered by the presence of obstacles or interference acquainted with a battlefield.

The subsequent versions of the system may look into adopting more sophisticated communication systems and hardened components to address this drawback. Then again, even with these, the preliminary outcomes are encouraging and suggest that this setup is poised to be a great asset in the field of military operations.

### E. *Model Performance and Evaluation Metrics*

The main measure applied for evaluation of the model was accuracy which refers to the number of audio samples which are identified rightly (cow, cow nutrition) whether authentic or forged. Additionally, other measures were applied such as precision, recall and f1 score in order to gauge more effectively the extent to which the model was able to carry out its objective, particularly voice forgery identification. In addition to this, we also performed a visualization of the performance of the model in terms of the confusion matrix where, the number of accurately identified genuine voices is plotted, and also the number of correctly identified fakes.

### F. *Feature Extraction*

In order to train our model to recognize whether a voice is real or AI-generated, we needed to break down the audio into key features that the model can learn from. The process of feature extraction is crucial because it transforms raw audio data into a form that is easier for the model to process and understand. For this, we used a technique known as Mel-frequency cepstral coefficient (MFCCs), which is one of the most widely used methods in speech processing. The reason MFCCs are so popular in speech recognition and speaker verification is that they mimic how the human ear perceives sound. Our ears are more sensitive to certain frequencies than others, and MFCCs are designed to emphasize these frequencies, making it easier to tell voices apart.

For every audio sample in our dataset, we utilized a tool called Librosa, a Python library for audio analysis, to extract 13 MFCCs. Here's how the process works in a simplified way:

**Breaking the audio into short frames:** In this step, the audio is cut into very short elements called frames. Each frame is typically shorter than a second, in most cases between 20 and 40 milliseconds. This is done since over time speech changes and therefore by segmenting it we are able to focus on the changes that occur in the voice for the duration of the given sample.

**Fourier Transform:** Fourier Transform is then performed on each of the frames. In this case, the audio is transformed from the time domain where we can only see the waveforms to the frequency domain where we can see the various frequencies of sound contained in the audio. This is an important step for analysing how high or low the voice of the speaker is. **Mel Scale:** The Mel scale is then applied when audio signals are transformed from time-domain to frequency-domain, and it transforms the inherent distribution of frequencies by compressing the high frequencies since human ears do not respond well to them. After applying the mel scale, the sound may be enriched however the focus of the model will steer towards the better visualized, understandable components of the sound. **Cepstrum Calculation:** Lastly, we obtain the cepstral distance which tells us how varying energies are spent at various frequencies of an audio signal. This process yields the basic parameters for speech, the 13 MFCCs associated with the speech signal. Moving on to the methodology employed for the analysis, the objective is to ascertain if the given voice has been recorded as real human voice or has been created by technology in the form of artificial intelligence. MFCCs in this case help a model to emphasize the elements of the voice of the speaker which will be very difficult for any deep fake to synthesized naturally. Because MFCCs contain both fine level of information like pitch, and information like voice patterns, how the voice changes over time, they equip the model with data that will enable it predict accurately e.g. deepfake voices.

### G. Data Management

After extraction of the MFCCs, the values were structured and stored into feature matrices such that each matrix comprised of number of rows corresponding to the number of video frames and the columns representing the 13 MFCCs respectively. Such matrices were prepared and forwarded to the LSTM network during its training as well as testing.

Employing MFCCs, we were able to shrink a vast array of complex audio data into one that is easily comprehensible to the model thereby giving it latitude on the regions of interest on the voice signal and not flooding it with irrelevant data.

## V. EXPERIMENTAL RESULTS

After training the LSTM-based speaker verification model for 50 epochs, we tracked both the training accuracy and validation accuracy to ensure the model was learning effectively and generalizing to unseen data. The following figure shows the accuracy and loss curves over the course of the training process.

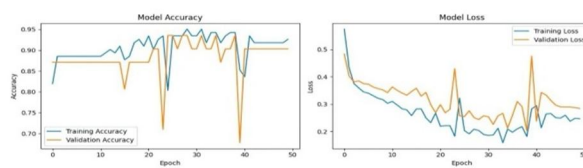


Figure 1: Training and validation Accuracy and Loss

The accuracy curve shows that the model steadily improved its performance during training, reaching a validation accuracy of around 85%. The loss curve also demonstrates a decrease in both training and validation loss, although some fluctuations were observed. This indicates that the model was effectively learning to classify voices as real or fake.

### A. Confusion Matrix

To better understand the model’s performance, we analyzed the confusion matrix, which illustrates the model’s predictions for both real and fake voices.

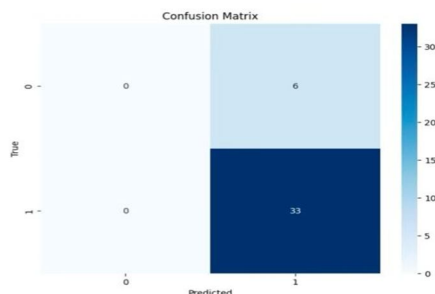


Figure 2: Confusion Matrix for the Test Dataset

From the confusion matrix, we can see that the model was highly effective at identifying fake voices (class 1) with perfect accuracy, but it struggled to correctly identify real voices (class 0), misclassifying all of them as fake. This issue could be addressed by further tuning the model or balancing the dataset more effectively.

### B. Classification Report

The classification report provides additional metrics such as precision, recall, and F1-score. These metrics help us better understand the trade-off between false positives and false negatives.

	0	0.00	0.00	0.00	6
	1	0.85	1.00	0.92	33
accuracy				0.85	39
macro avg	0.42	0.50	0.46		39
weighted avg	0.72	0.85	0.78		39

Figure 3: Precision, Recall, and F1-Score

The model achieved an overall accuracy of 85%. However, the performance for real voices was less effective, with precision and recall both being 0. This suggests that while the model performs well on fake voices, it needs improvement in accurately identifying real ones. The F1-score for fake voices was 0.92, showing that the model can reliably detect AI-generated speech.

### C. Preprocessing

Before feeding the audio data into the model, it was necessary to preprocess the raw audio files to extract meaningful features that the model could learn from. Proper preprocessing is a critical step in ensuring that the model can accurately distinguish between real and deepfake voices.

This preprocessing pipeline ensured that the audio data was in an optimal format for training the LSTM model. The careful extraction of MFCCs, padding of sequences, and application of data augmentation techniques provided the model with rich, diverse input data, helping to improve its performance in detecting both real and AI-generated voices.

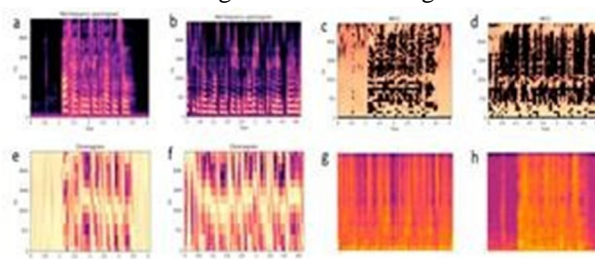


Fig 3. (a) original Mid-freq audio, (b) cloned Mid-freq audio, (c) original MFCC audio, (d) cloned MFCC audio, (e) original chromagram audio, (f) cloned chromagram audio, (g) original spectrogram audio, (h) cloned spectrogram audio.

## VI. FUTURE SCOPE

- 1) *Latency Minimization:* In order to achieve the expected level of performance in real time with little or no latency, the future systems may use faster communication standards such as 5G or LoRa along with faster optimized algorithms.
- 2) *Wearable Technologies:* To augment situational awareness within the battlefield, wearables like health monitors and visual cue-based augmented reality can be added to the system.
- 3) *Machine Learning:* Gesture sets can be trained per user using machine learning which will recognize their gestures and hence, enhancing the flexibility and accuracy of the system with time.
- 4) *Advanced Sensors:* The use of multi-modal sensors like electromyography (EMG) sensors can also enhance the performance of the system in detecting the various gestures. The system can also be developed to be adaptive in harsh environmental conditions, such as those found in the modern battlefield.
- 5) *Security:* Future system editions should include strict security standards to protect sensitive military information.
- 6) *Power Management:* On the other hand, energy efficiency can be focused on developing energy-harvesting techniques, like the use of body heat, kinetic energy, and so on. And also to enhance the machine's work time, Low-powered elements may be installed.
- 7) *Extended Applications:* There is high potential for further modification of the system to address other applications such as in wound care telemedicine in combat scenarios, remote operation of drones and espionage.



## VII. FUTURE SCOPE

The development of advanced speaker verification models has significant potential to impact mankind and enhance everyday life in several meaningful ways:

### A. *Enhanced Security and Privacy*

As deepfake technology continues to evolve, the ability to distinguish between real and AI-generated voices will become crucial for protecting people's privacy and ensuring secure communication. Future versions of this model could be integrated into everyday applications such as voice-activated assistants, banking systems, and smart home devices, providing an extra layer of protection against identity theft and fraud.

### B. *Trust in Digital Communication*

In a world where digital interactions are becoming more common, ensuring that the person speaking on the other end is who they claim to be is essential. This technology could play a vital role in improving the trustworthiness of voice communications, especially in sensitive contexts like remote work, telemedicine, and legal proceedings.

### C. *Accessible Authentication for All*

Voice authentication can make secure access easier and more inclusive, especially for individuals with physical disabilities or those who have difficulty using traditional methods like passwords or biometric scans. By improving the detection of fake voices, future systems could make everyday tasks like unlocking a phone, authorizing payments, or accessing personal accounts more secure and convenient.

### D. *Protecting Public Figures and Content Creators*

With the increasing use of voice cloning for malicious purposes, such as impersonating public figures or content creators, this technology could help safeguard their reputations and public trust. For example, social media platforms and news outlets could adopt such systems to verify that the voice used in a clip belongs to the actual speaker, reducing the spread of misinformation.

### E. *Safe Use of AI in Entertainment and Media*

As AI voice generation is becoming more prevalent in media and entertainment, there is a need to ensure that it is used ethically and safely. This model could be further developed to help monitor and regulate AI-generated content, ensuring that voices are not used without consent or in misleading ways.

### F. *Future-Proofing Against Technological Threats*

As deepfake and other AI technologies continue to advance, systems like this will be critical in future-proofing security measures. By staying ahead of these technologies, the development of more sophisticated speaker verification models could protect both individuals and organizations from the misuse of AI in voice manipulation.

## VIII. CONCLUSION

In this study, we developed a speaker verification model using LSTM networks and MFCCs to distinguish between real and AI-generated voices. The model achieved an accuracy of 84.62%, showing strong performance in detecting fake voices with an F1-score of 0.92. However, it struggled with identifying real voices, indicating a need for further optimization.

By applying data augmentation techniques, we improved the model's robustness, making it better suited for real-world applications. This research lays a foundation for enhancing voice authentication systems, and future work could focus on improving real voice classification and exploring more advanced architectures.

## IX. ACKNOWLEDGMENT

We express our deepest appreciation to Dr. Supriya Telsang for her invaluable guidance and support during the course of the "Speaker Verification Model using LSTM" project. Her knowledge and motivation were key in fulfilling this project. We also thank the student testers for their feedback which assisted in enhancing the platform. We also thank our colleagues and mentors who supported us at different stages of this project. Finally, we value the commitment and efforts of the entire development team that made this project a reality.



10.22214/IJRASET



45.98



IMPACT FACTOR:  
7.129



IMPACT FACTOR:  
7.429



# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24\*7 Support on Whatsapp)