

# Audio Impersonation detection through Spectrogram Insights

**MM. Nirmala**, Asst. Professor , Kallam Haranadhareddy Institute of Technology,  
Chowdavaram, Guntur, Andhra Pradesh, India-522019

**KUNAPAREDDY YASASWINI DEVI**-[yasaswinidevikunapareddy@gmail.com](mailto:yasaswinidevikunapareddy@gmail.com),  
**NISSANKARARAO VENKATA TARUN, GOLI VENKATA RAMANJI, TADIBOINA  
VENKATA SURENDRA GOPAL**, Department of Information Technology, Kallam  
Haranadhareddy Institute of Technology, Chowdavaram, Guntur, Andhra Pradesh, India-522019

**Abstract:** The proliferation of deepfake audio and other synthetic audio technology has made it very difficult to determine whether audio recordings are genuine. In light of these difficulties, this study seeks to use spectrogram analysis in conjunction with a Support Vector Machine (SVM) classifier to identify synthetic sounds. The technology transforms audio files into spectrograms, which record the audio's frequency content as it evolves, in order to process them. The SVM model is trained using these spectrogram features, and it then determines if the audio is real or not.

**Index terms** - *Audio Classification, Frequency Content, Spectrogram Features, Support Vector Machine (SVM), Audio Authentication, Deepfake Audio, and Synthetic Audio.*

## 1. INTRODUCTION

Rapid advancements in deepfake technology in the last several years have made it possible to produce footage that is both incredibly lifelike and heavily edited [4]. Although deepfakes involving the alteration or replacement of a person's face in films have received a lot of attention, audio deepfakes are just as dangerous [5]. Serious consequences, including disinformation, identity theft, financial scams, and cyber dangers, can result from these deepfakes since they artificially produce or alter speech in a way that convincingly mimics an individual's voice [4]. One of the biggest problems with digital communication is the ease with which people can alter audio recordings, which makes it hard to tell if what you're listening to is real [2]. Facial recognition, motion inconsistency analysis, and anomaly detection are some of the visual analytic techniques that have traditionally been utilized in deepfake detection [4]. Despite the effectiveness of these methods in identifying edited movies, they frequently miss tiny changes in audio content [5]. But there is hope in audio deepfake detection, which looks at the inherent features of human voices, such as speech patterns and frequency changes, to identify fake voices [1]. It is possible to detect instances of

artificial generation or manipulation of an audio sample by analyzing these features [5].

An audio deepfake detection system based on Mel Frequency Cepstral Coefficients (MFCCs) and a Support Vector Machine (SVM) classifier with a Radial Basis Function (RBF) kernel is proposed in this research to tackle this difficulty [7]. For the purpose of identifying irregularities in synthetic audio, MFCCs find extensive application in speech processing due to their ability to accurately capture the fundamental spectrum and frequency characteristics of human speech [1]. Also, by using derived feature patterns to successfully differentiate between real and fraudulent audio samples, SVM with the RBF kernel provides strong classification capabilities [7].

Principal Component Analysis (PCA) is also applied to the retrieved features to decrease computational complexity and increase model efficiency [3]. By assisting with dimensionality reduction, principal component analysis (PCA) enables the model to zero in on the most important data points while removing unnecessary ones [3]. As a result, the classification process is more efficient and accurate [8]. Creating a trustworthy system that can distinguish between authentic and deepfake audio recordings is the primary objective of this research [5]. An effective response to the growing threats of synthetic voice impersonation and modified speech content can be achieved by utilizing powerful machine learning algorithms and audio feature analysis. This system seeks to do just that [6].

In domains including news broadcasting, fraud detection, and cybersecurity, our work adds to the continuing efforts to strengthen digital media security and avoid the abuse of deepfake technologies [4].

## 2. LITERATURE SURVEY

[1] The basic ideas of digital voice processing are discussed in this book by Rabiner and Schafer (2011). Speech synthesis and recognition are just two of the many fields covered, along with methods for representing and processing speech signals. Additional

topics covered in the book include feature extraction for voice analysis, digital filtering, and spectral analysis. Researchers and engineers involved in voice processing will find it to be an exhaustive resource.

[2] Natural language processing (NLP), speech recognition, and other areas of language and speech processing are thoroughly covered in Jurafsky and Martin's (2021) comprehensive introduction. For text and audio analysis, the book delves into linguistic techniques, deep learning, and machine learning. It goes on to cover statistical modeling, voice signal processing, and its practical uses in translation and virtual assistants. Many computational linguists and AI researchers look to the book as a starting point for their work.

Principal Component Analysis (PCA) is a popular statistical method for reducing dimensionality, as described by Jolliffe (2002). In this book, the author delves into the theory behind principal component analysis (PCA), as well as its many real-world applications. Use of principal component analysis (PCA) for feature selection, data visualization, and noise reduction is described. In addition to covering PCA and its uses in machine learning, the book delves into its expansions.

[4] In their review of deepfake audio detection approaches, Wu and Lyu (2022) bring attention to both the difficulties and the recent successes in this field. Synthetic or altered speech can be detected using machine learning and deep learning algorithms, which are reviewed in this study. Issues with adversarial detection, assessment metrics, and datasets are also covered. A road map for future research on deepfake audio threats is provided by the paper.

Methods for identifying speech deepfakes that rely on machine learning are suggested by Patel and Patel (2020). Their research compares various models and characteristics for accurate speech-to-text conversion. Topics covered in the article include classifier efficiency, feature engineering, and spectrogram analysis. The significance of diverse and strong datasets in deepfake detection is further emphasized.

[6] The Python package Librosa is introduced by McFee et al. (2015) for the purpose of analyzing audio and music signals. Machine learning feature extraction, visualization, and processing of signals are some of its functions that are covered in the study. It gives a general outline of methods for filtering, spectral modifications, and time-series representation. Research on music information retrieval and speech processing makes heavy use of Librosa.

[7] Support Vector Machines (SVMs) are a potent supervised learning technique that Cortes and Vapnik (1995) introduce. This article delves into the mathematical construction of support vector machines (SVMs) and how they perform on categorization jobs. Kernel functions, optimization methods, and maximizing margins are highlighted. Pattern recognition, audio processing, and bioinformatics have all seen extensive use of support vector machines (SVMs).

A popular Python package for machine learning, Scikit-learn was introduced by Pedregosa et al. (2011). Classification, regression, clustering, and model evaluation methods are some of its features that are covered in the study. It highlights the scalability, simplicity, and integration of the library with SciPy and NumPy. Data science scholars and practitioners now rely on scikit-learn.

Streamlit is a free and open-source Python framework for creating dynamic web apps; its documentation is available at [9] Streamlit (2024). It provides a user-friendly interface for deploying data visualizations and machine learning algorithms. Methods for optimizing efficiency, API references, and detailed instructions are all part of the documentation. Data scientists and AI developers rely on Streamlit for quick prototyping. [10] A tutorial on how to install Streamlit apps on Streamlit Cloud may be found on GitHub (2024). Learn all you need to know about setting up and configuring Streamlit programs to run online with our comprehensive guide. It delves into authentication, scalability, and deployment pipelines. With the help of this manual, programmers can set up and manage their cloud-based Streamlit projects with ease.

### 3. METHODOLOGY

#### i) Proposed Work:

Integration of state-of-the-art feature selection and classification algorithms improves deepfake audio detection in the suggested system. To improve computing performance while keeping crucial information, Principal Component Analysis (PCA) [3] is used to reduce dimensionality.

Additionally, the system makes use of Support Vector Machines (SVM) with the Radial Basis Function (RBF) kernel [7], which allows for better categorization by efficiently handling complex, non-linear data distributions. To further capture the spectral and temporal aspects of speech, it uses Spectrogram-based features in conjunction with Mel-Frequency Cepstral Coefficients (MFCCs) [1, 6]. When compared to traditional models that depend only on less complex classification algorithms and feature extraction approaches, this multi-feature approach

greatly improves detection accuracy, making the system more robust [5].

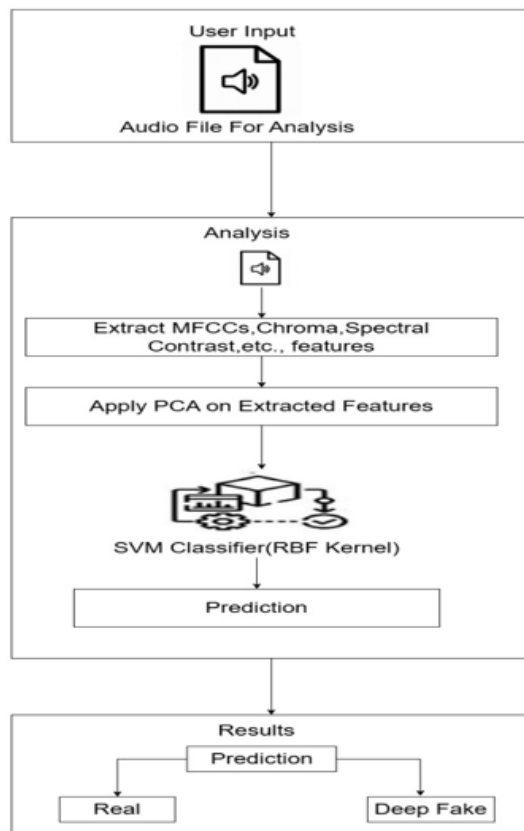


Fig 1 Proposed Architecture

- The suggested technique extracts features based on spectrograms and MFCCs, which capture the time and frequency variations in speech more precisely, allowing for a more accurate representation of the features. This improves the accuracy of spotting impersonation attempts by detecting subtle deepfake artifacts that MFCCs alone could overlook.
- The RBF Kernel for SVM Improves Classification: The RBF kernel records nuanced correlations between actual and false voices, in contrast to Linear SVM, which has trouble with complicated, non-linear data. This improves the model's classification performance by enabling it to efficiently distinguish between various forms of deepfake speech.
- Efficient and Quicker with Principal Component Analysis (PCA): By eliminating unnecessary and redundant characteristics, PCA makes the system more efficient and quick without sacrificing valuable data. Not only does this avoid overfitting, but it also speeds up testing and training.

- Hyperparameter tweaking (C & Gamma for SVM) guarantees higher dependability, leading to stronger performance evaluation and model tuning. Its efficacy can be better gauged with the use of evaluation metrics like as Accuracy, Precision, Recall, and F1-Score.

The user's interaction with the system through a front-end interface initiates the project workflow. The user is prompted to upload the model file before proceeding unless there is already a trained model available; otherwise, they can just upload an audio file for analysis. Standardizing the format, converting to mono, and adjusting the sampling rate are all steps in the preprocessing of the uploaded audio that are done to make sure it is consistent.

Then, features based on spectrograms and time-domain properties, such as MFCCs, are extracted. Applying Principal Component Analysis (PCA) [3] improves computational efficiency and model performance by eliminating redundancy. After that, the cleaned-up set of features is fed into a Support Vector Machine (SVM) using an RBF kernel [7], which improves classification accuracy by capturing complicated data relationships.

Using the features that were retrieved, the trained model determines if the audio is real or not. Results are shown on the front-end interface by the system, with clarity provided by visualizations such as spectrogram plots, confusion matrices, or probability scores. Feature selection and advanced classification algorithms are integrated in this approach to guarantee accurate and robust deepfake audio detection [1, 5].

#### ii) System architecture:

Users can upload an audio file for testing through the frontend, which is the Streamlit web app.

- The system has a client-server architecture.
- To process and categorize the audio, the backend uses a Python server with an ML model.
- You have the choice to store audio samples and findings in the database or storage.

From extracting features to deploying models, this chapter details the entire process of implementing the audio impersonation detection system.

#### Feature Extraction & Preprocessing:

It is necessary to transform unprocessed audio data into features before training a model.

##### Feature Extraction:

Librosa, an audio analysis library for Python, is used to do feature extraction. Using the extracted features, we can tell the difference between real and fraudulent audio samples.

##### Preprocessing Steps

- First, obtain the dataset, which contains both real and synthetic audio samples.

- The second step is to use the `extract_features()` function to extract features.
- Third, save the input features (the features that were extracted) in X and the labels in y.

Principal Component Analysis (PCA) is used to reduce dimensions and remove redundancy while maintaining critical information. This is done since the extracted feature set is high-dimensional.

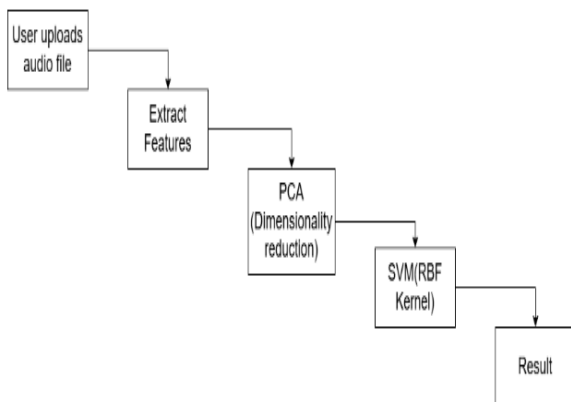
Using PCA

- Makes training more efficient by reducing feature complexity.
- Enhances model performance.
- Prevents overfitting by eliminating noise.

#### Classification Algorithm:

Machine Learning Support Vector (SVM) using RBF Kernel Its proficiency in dealing with data with many dimensions led to the selection of the Support Vector Machine (SVM) classifier.

**Model Training & Optimization:** The SVM model is fine-tuned by means of hyperparameter tweaking with Randomized Search CV.



Accuracy, precision, recall, and F1-score are used to assess the model's performance after training.

**Accuracy:** The reliability of a test is defined as its ability to differentiate between healthy and sick individuals. Find the percentage of cases that were true positives and true negatives to get an idea of the test's accuracy. Formally speaking, this is:

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN} \quad (1)$$

**Precision:** A high level of precision indicates that the majority of instances or samples were correctly classified. The precision is determined by applying the formula:

$$Precision = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}} \quad (2)$$

**Recall:** Recall is a measure of a model's ability to identify all instances of a class that are relevant in

machine learning. If the ratio of correctly predicted positive observations to total positives is high, then the model adequately captures instances of the class.

$$Recall = \frac{TP}{TP + FN} \quad (3)$$

**F1-Score:** F1 score is a metric for evaluating the accuracy of machine learning models. Integrating the scores for model recall and precision. A model's accuracy is defined as the percentage of times its predictions were correct over the whole dataset.

$$F1\ Score = 2 * \frac{Recall * Precision}{Recall + Precision} * 100 \quad (4)$$

#### Deployment Plan:

Sending out The trained model is prepared for use in the actual world through the deployment phase. Here, we construct an interactive web tool for detecting audio impersonation using Streamlit, a Python-based framework.

#### Algorithm:

**Support Vector Machine (SVM) with RBF Kernel** Since it excels at dealing with data with many dimensions, the Support Vector Machine (SVM) classifier was used. The complicated patterns found in thyroid illness data are well-suited for classification using Support Vector Machine, because to its success in high-dimensional areas [20]. Improving diagnostic prediction accuracy is the goal of finding the best hyperplane that divides classes. One possible form of the equation representing the linear hyperplane is:

$$wTx + b = 0 \quad (5)$$

#### RBF Kernel Use

- Manages data that is not separable in a linear fashion. Optimal decision boundaries are found by mapping input features onto a higher-dimensional space.

## 4. EXPERIMENTAL RESULTS

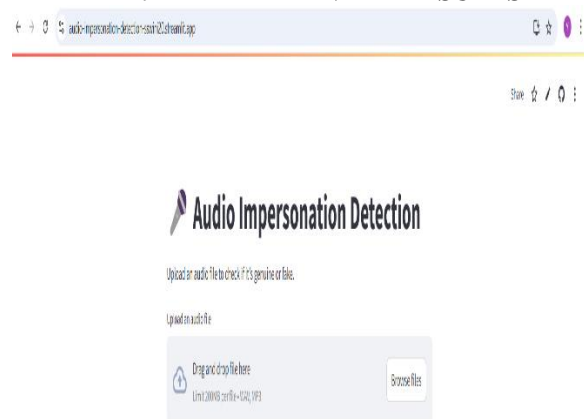


Fig 2 Home page



## Audio Impersonation Detection

Upload an audio file to check if it's genuine or fake.

Upload an audio file

Drag and drop file here  
Limit 200MB per file • WAV, MP3

Browse files

real 203.mp3 121.3KB

X

0:00 / 0:10

Prediction: ● Genuine

Confidence: 98.82%

Fig 3 Prediction result genuine

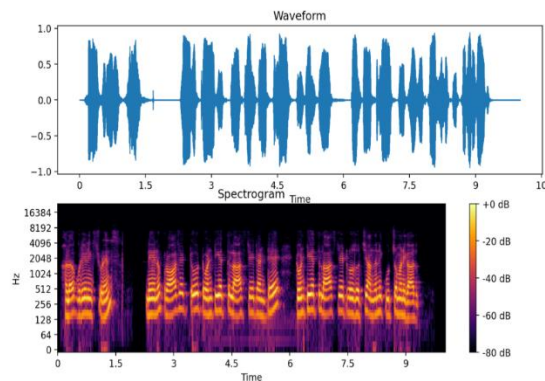


Fig 4 Waveform and Spectrogram Representation of an Audio Signal-1

## Audio Impersonation Detection

Upload an audio file to check if it's genuine or fake.

Upload an audio file

Drag and drop file here  
Limit 200MB per file • WAV, MP3

Browse files

VID-20250316-WAO015.wav 0.6MB

X

0:00 / 0:03

Prediction: ● Fake

Confidence: 77.91%

Fig 5 Prediction result fake

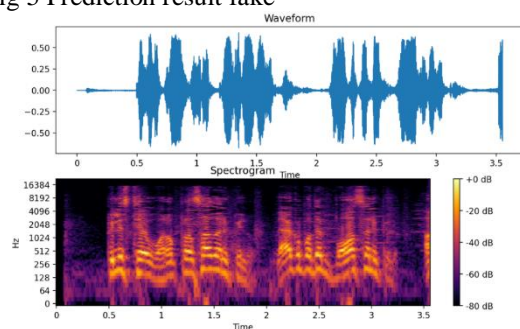


Fig 6 Waveform and Spectrogram Representation of an Audio Signal-2

## 5. CONCLUSION

Using machine learning approaches, the audio impersonation detection project determines if a voice recording is authentic or not. The system guarantees accurate detection by using MFCC to extract features, PCA to reduce dimensionality, and SVM with RBF kernel for classification. Short audio samples may be handled successfully by the project, which helps to avoid misclassification. The model is even more user-friendly when launched as a web app; users can simply input audio files and get immediate results. This research helps to improve the security of voice-based authentication systems, which has applications in areas such as fraud detection, cybercrime investigations, deepfake prevention, and voice authentication security. Its potential for future enhancements like real-time detection, multi-language support, and the creation of mobile applications is immense. The project's realistic and effective method of identifying voice impersonation makes it an asset to AI-driven speech processing and security systems.

## 6. FUTURE SCOPE

Possible upgrades to the Audio Impersonation Detection System in the future include:

1. Integrating Deep Learning: Improving audio feature extraction and classification with CNNs or LSTMs.
2. The second feature is real-time recording, which means users can record their voice without having to submit a file by using the web app itself.
3. Using SVM in conjunction with additional classifiers, such as Random Forest or XGBoost, to achieve higher accuracy is the third method, multi-model classification.

## REFERENCES

1. Rabiner, L. R., & Schafer, R. W. (2011). Theory and Applications of Digital Speech Processing. Pearson.
2. Jurafsky, D., & Martin, J. H. (2021). Speech and Language Processing (3rd ed.). Pearson.
3. Jolliffe, I. T. (2002). Principal Component Analysis. Springer.
4. Wu, Y., & Lyu, S. (2022). Deepfake Audio Detection: A Survey. arXiv preprint arXiv:2201.07432.
5. Patel, H., & Patel, D. (2020). Voice Deepfake Detection using Machine Learning. International Journal of Computer Applications, 175(7), 23-29.
6. McFee, B., et al. (2015). Librosa: Audio and Music Signal Analysis in Python. Proceedings of the 14th Python in Science

- Conference. Retrieved from  
<https://librosa.org/doc/latest/>
7. Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine Learning*, 20(3), 273-297.
  8. Pedregosa, F., et al. (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12, 2825-2830. Retrieved from <https://scikit-learn.org/stable/>
  9. Streamlit. (2024). Streamlit Documentation. Retrieved from <https://docs.streamlit.io/>
  10. GitHub. (2024). Streamlit Cloud Deployment Guide. Retrieved from <https://docs.streamlit.io/streamlit-cloud>