



# Deepfake Audio Detection Using Ensemble CNN-BiLSTM Architecture

<sup>1</sup>Dr. J Sharmila Joseph, <sup>2</sup>Anaswar L, <sup>3</sup>Harshwardhan Hande and <sup>4</sup>Soham Dange

<sup>1</sup>Professor, AIML, VIT Bhopal, Vellore Institute of Technology, Madhya Pradesh, India.

<sup>2,3,4</sup>Student, CSE SCOPE, VIT Bhopal, Vellore Institute of Technology, Madhya Pradesh, India.

## Abstract

A new approach uses deep learning to spot fake audio made by artificial intelligence. Instead of one single network, it combines convolutional neural nets with bidirectional long short-term memory models working together. Features pulled from sound include things like mel-frequency cepstral coefficients, mel-spectrograms, chroma data, and how pitch changes over time. While the CNN looks at mel-spectrograms to find fixed visual-like patterns in the audio, the Bi-LSTM focuses on sequences of MFCCs. Temporal quirks common in computer-made voices are caught through this sequence analysis. Together, these parts help tell real speech apart from synthetic versions.

**Keywords:** Multimodal Feature Fusion, Hybrid Deep Learning, Audio Anti-Spoofing Countermeasures.

## 1. Introduction

New progress in turning text into speech and changing voices digitally now lets machines sound like real people. Because of this shift, fake spoken content is harder to spot than before. These tools can be helpful in everyday tasks yet open doors to misuse at the same time. Problems show up when systems trust a voice as proof of identity. Media files might carry false messages without clear signs. Security online faces new hurdles because of altered audio tracks. Telling real recordings from fakes has grown more urgent by the day. Researchers now focus energy on spotting clues hidden in sound patterns. The rise of realistic machine-made voices changes how we check what's true.

Back then, spotting fake voices meant using preset sound traits - things like MFCCs, LPC values, because those captured how speech looked in frequency patterns [9, 11, 34]. Even if they worked somewhat well picking up traces left by older voice generators, once newer systems arrived, results got shaky [4, 22]. Because of that shortcoming, scientists lately shifted toward neural networks instead; now models learn meaningful clues straight from raw audio or time-frequency visuals without being told what to look for [2, 24, 30].

This study presents a system built on combined deep learning methods, linking CNNs with Bi-LSTM layers to boost accuracy in spotting fake audio. Spectrogram visuals are processed effectively by CNNs, which detect fine-grained frequency shifts across sound inputs [2]. Temporal sequences, meanwhile, gain sharper representation through recurrence-based nets - LSTMs and their bidirectional variants handle timing links in data streams [32, 33].

What stands out in this work comes down to a few key points  
A split-path design pulls apart how sound changes over time

\*Corresponding Author: Anaswar L

from its frequency makeup. One path tracks rhythms as they unfold moment by moment. The other locks onto pitch patterns across different bands. Together, they build a fuller picture without merging too soon. Each branch works on its own rhythm. Information flows side by side rather than stacking early. This separation helps preserve subtle differences in voice quality.

A setup where a CNN picks up visual patterns while a Bi-LSTM tracks time-driven changes across frames - borrowing ideas from blended deepfake spotting systems [12, 14].

Most decisions come from combining several model outputs, boosting accuracy by relying on group results instead of single guesses - this mirrors methods seen in anti-spoof studies [5].

## 2. Literature Review

One reason progress in speech synthesis matters? Tools like TTS and VC now make real and fake voices hard to tell apart. Because of that, new ways to spot fakes have started appearing - aimed at strengthening speaker checks. Some key steps forward look like this:

### i). Foundational Benchmarks and Datasets

Starting off, progress in catching fake audio got a big boost from the ASVspoof challenges. These events set up common rules plus shared data so researchers can test their systems fairly. One key resource? The ASVspoof 2019 dataset [1], now standard across many studies. It packs thousands of faked voices - some stitched together, others mimicked or recorded then played back. Later on, the 2021 version [6] turned things up: tougher sound environments showed how well detectors handle real-world chaos. Not long before that, De Leon and team [10] dug into subtle clues hidden in voice patterns, spotting

gaps between real talk and computer-made copies.

### ii). Feature Engineering and Acoustic Robustness

Few people looked beyond manual methods when studying flaws in computer-made voices. Those approaches ruled the scene until smart algorithms took over. Crafted traits once guided every analysis - carefully shaped by human effort. Only later did pattern-learning systems begin replacing them. The shift started slowly, then accelerated without warning.

**Spectral Resolution:** Starting off differently, Todisco and team introduced something called Constant Q Cepstral Coefficients - these adjust how finely frequencies are measured. That shift helps catch odd sounds in audio that usually hint at fake speech. Meanwhile, work led by Kinnunen took a look at various ways to study sound patterns when spotting voice conversion tricks.

**Environmental Robustness** Starting off differently each time helps keep things fresh. Performance often drops when recordings come from varied setups or gadgets. Zhang's team tackled device differences by designing traits less swayed by hardware shifts. Their approach softens the effect of inconsistent capture tools. On another path, Wang's research looked at messy audio settings where fakes hide in static. They crafted techniques pulling fake signs apart from surrounding sound clutter.

### iii). Neural Network Architectures

Now machines spot fake audio better because deep learning digs into voice details without handcrafted rules. Hidden layers catch subtle clues regular programs miss. Training on raw sound lets the system build its own sense of what's real. Progress here means fewer errors when telling natural speech from synthetic. Earlier attempts struggled with variation, but today's networks adapt through exposure. Each pass through data sharpens the model's instincts. Real gains come not from more features, but deeper processing. Signals once too messy now reveal their truth under scrutiny.

**Convolutional Models:** Spectrogram visuals get processed by CNNs much like picture data does. Local patterns tied to fake voice signals were caught well through these networks, according to work by Sahidullah *et al.* [2]. Lighter versions came later - Michelsanti *et al.* [19] built simpler designs that kept accuracy high without needing heavy computing power.

**Temporal and Hybrid Models:** Speech timing helps spot odd patterns when someone talks. Wu's team looked into RNNs to track how fake voices change over time. After that, Korshunov's group tested a mix of CNN and LSTM setups - pulling out details across space while following shifts step by step through sequences.

**Focus in Transformer Networks:** Newer studies use systems

that help models pay closer attention to key parts of sound signals [15]. Because they handle patterns across time well, transformers now show up in fake audio spotting work [24]. Specific network designs tuned to frequencies aim at areas where fake traces often hide [23].

### iv). End-to-End Processing and Raw Waveforms

Starting fresh, some researchers skip preset sound details entirely - training systems straight from unprocessed audio signals. A different path emerges when raw waveform data shapes the learning, bypassing traditional feature extraction steps altogether.

Rather than relying on processed features, Jung and team introduced RawNet2 - a model trained straight from unedited sound waves. Meanwhile, work by Zhang's group dug into fake speech spotting using waveforms too, showing how learning everything at once beats older methods if catching tiny traces left by synthetic voice tools is the goal.

### v). Advanced Learning Strategies and Generalization

When deepfake techniques get smarter, making detectors work better across different cases grabs more attention in studies. While new tricks pop up, scientists spend extra time on tools that handle surprises well.

**Generalization Challenges**

What happens when new spoofing tricks show up? Aloufi *et al.* [4], along with Todisco *et al.* [22], pointed out how tough it is to catch these unknown types. This problem often goes by another name: the generalization gap.

**Learning Paradigms:** Chen and team [18] looked into how models can pick up patterns by themselves, skipping the need for tons of tagged examples. Another shift happens when knowledge moves across tasks - Kumar *et al.* [21] tested ways to repurpose models built on one set of data so they work elsewhere. **Fusion and Explainability**

Combining multiple detection strategies has proven effective in improving performance. Multi-feature fusion techniques [14] and ensemble learning approaches [5] enable models to integrate complementary information from different feature representations. At the same time, explainable AI techniques have been explored to ensure that detection decisions are based on meaningful spoofing artifacts rather than dataset bias [25].

## 3. Dataset Used

Dataset was acquired from the following link (<https://www.kaggle.com/datasets/adarshsingh0903/audio-deepfake-detection-dataset?resource=download&select=FlashSpeech>) which consisted of multiple real and fake audios in order to train and perform the model.

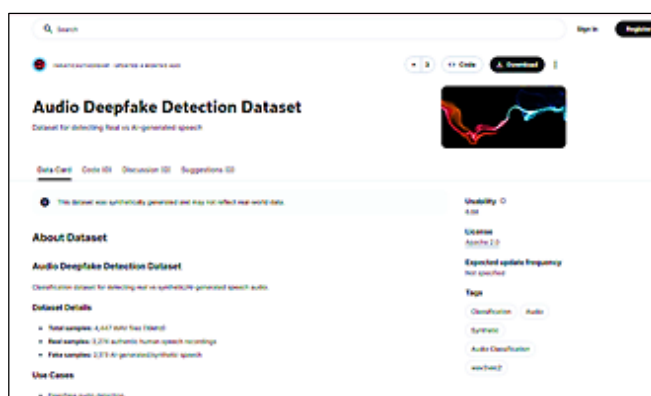


Fig 3.1: Dataset Page

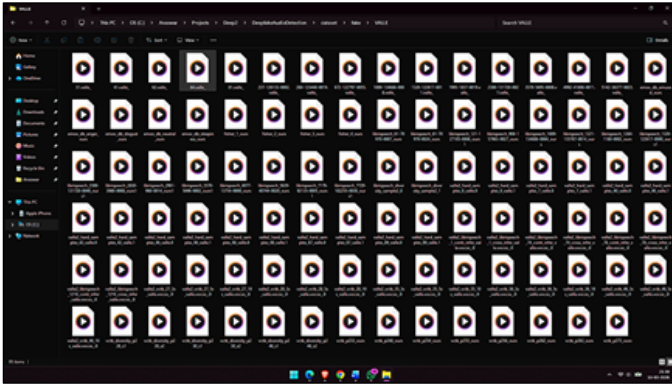


Fig 3.2: Valle – Fake Dataset

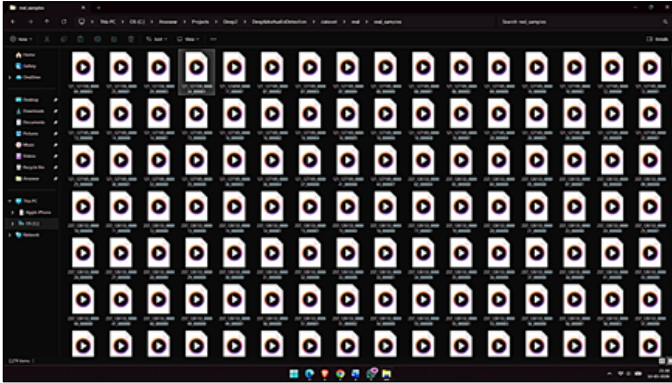


Fig 3.3: Real samples

#### 4. Proposed Workflow

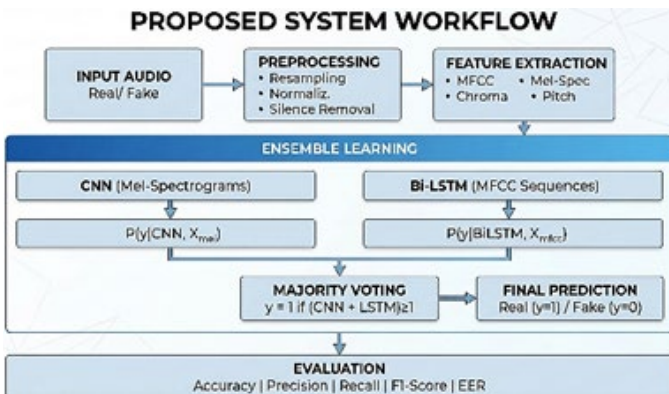


Fig 4.1: Proposed System Workflow

### 5. Architecture Diagram

#### 5.1. System Architecture

The system architecture diagram of the proposed model is shown below.

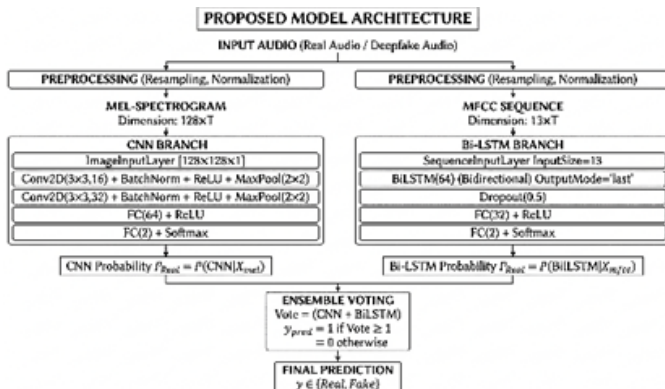


Fig 5.1: Model Architecture

#### 5.2. CNN Architecture Details

Table 5.1: CNN Architecture Details

Layer	Type	Parameters	Output Size
1	Input	-	128×128×1
2	Conv2D + BatchNorm + ReLU	16 filters, 3×3	128×128×16
3	MaxPooling	2×2, stride 2	64×64×16
4	Conv2D + BatchNorm + ReLU	32 filters, 3×3	64×64×32
5	MaxPooling	2×2, stride 2	32×32×32
6	FullyConnected + ReLU	64 units	64
7	FullyConnected + Softmax	2 units	2

#### 5.3. Bi-LSTM Architecture Details

Table 5.2: Bi-LSTM Architecture Details

Layer	Type	Parameters	Output Size
1	Input	Sequence (13 features)	13×T
2	BiLSTM	64 hidden units	128
3	Dropout	0.5	128
4	FullyConnected + ReLU	32 units	32
5	FullyConnected + Softmax	2 units	2

### 6. Mathematical Formulation

#### 6.1. Feature Extraction

##### 6.1.1. Mel-Spectrogram

The mel-spectrogram is calculated by applying mel-scaled filterbanks to the Short-Time Fourier Transform (STFT) of the audio signal:

$$S_{\text{mel}}(m, n) = \sum_{k=0}^{N-1} |X(k, n)|^2 \cdot H_m(k) \quad \text{Equation 6.1}$$

Where:

- $X(k, n)$  is the STFT of the audio signal at frequency bin  $k$  and time frame  $n$
- $H_m(k)$  is the  $m$ -th mel filterbank
- $N$  is the FFT size

The mel frequency scaling is given by:

$$f_{\text{mel}} = 2595 \cdot \log_{10} \left( 1 + \frac{f}{700} \right)$$

Equation 6.2

$$f = 700 \cdot (10^{f_{\text{mel}}/2595} - 1)$$

Equation 6.3

##### 6.1.2 MFCC Computation

MFCCs are computed by applying the Discrete Cosine Transform (DCT) to the log-mel spectrogram:

$$c[n] = \sum_{k=1}^K \tilde{s}[k] \cos \left( \frac{\pi n(k-0.5)}{K} \right), n = 1, 2, \dots, N$$

Equation 6.4

Where:

- $c[n]$  is the  $n$ -th MFCC coefficient

- $\tilde{s}[k] = \log(S_{\text{mel}}(k))$  is the log-mel spectrogram
- $K$  is the number of mel filterbanks
- $N$  is the number of desired cepstral coefficients (typically 13)

## 6.2. CNN Model

### 6.2.1. Convolutional Operation

The convolution operation in layer  $l$  is defined as:

$$y_{i,j}^{(l)} = \sigma \left( \sum_m \sum_n w_{m,n}^{(l)} \cdot x_{i+m,j+n}^{(l-1)} + b^{(l)} \right)$$

Equation 6.5

Where:

- $w_{m,n}^{(l)}$  is the kernel weight at position  $(m, n)$  in layer  $l$
- $x^{(l-1)}$  is the input feature map
- $b^{(l)}$  is the bias term
- $\sigma(\cdot)$  is the ReLU activation:  $\sigma(z) = \max(0, z)$

### 6.2.2. Batch Normalization

Batch normalization normalizes the activations:

$$\hat{z} = \gamma \cdot \frac{z - \mu}{\sqrt{\sigma^2 + \epsilon}} + \beta$$

Equation 6.6

Where:

- $z$  is the input (pre-activation)
- $\mu$  is the batch mean
- $\sigma^2$  is the batch variance
- $\gamma, \beta$  are learnable scale and shift parameters
- $\epsilon = 10^{-8}$  is a small constant for numerical stability

## 6.3. Bi-LSTM Model

### 6.3.1. LSTM Cell Equations

The LSTM cell computes the following operations at time step  $t$ :

Forget Gate:  $f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$

Input Gate:  $i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$

Cell State Update:  $C_t = f_t \cdot C_{t-1} + i_t \cdot \tilde{C}_t$

Output Gate:  $o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$

Hidden State:  $h_t = o_t \cdot \tanh(C_t)$

Where:

- $\sigma(\cdot)$  is the sigmoid function:  $\sigma(z) = \frac{1}{1+e^{-z}}$
- $W_f, W_i, W_c, W_o$  are weight matrices
- $b_f, b_i, b_c, b_o$  are bias vectors
- $[\cdot, \cdot]$  denotes concatenation

### 6.3.2. Bidirectional Extension

The bidirectional LSTM processes sequences in both forward and backward directions:

$$\vec{h}_t = \text{LSTM}(x_t, \vec{h}_{t-1})$$

$$\overleftarrow{h}_t = \text{LSTM}(x_t, \overleftarrow{h}_{t+1})$$

$$h_t = [\vec{h}_t; \overleftarrow{h}_t]$$

Equation 6.7

The final hidden state  $h_t$  concatenates both forward and backward hidden states, capturing context from both directions.

## 6.4. Ensemble Voting

### 6.4.1. Classification Probability

Let  $P_{\text{CNN}}(y = 1 | X_{\text{mel}})$  be the probability predicted by the CNN model given mel-spectrogram input  $X_{\text{mel}}$ , and  $P_{\text{BiLSTM}}(y = 1 | X_{\text{mfcc}})$  be the probability from the Bi-LSTM model given MFCC sequence  $X_{\text{mfcc}}$ .

### 6.4.2. Majority Voting

The final prediction is determined by majority voting:

$$\text{Vote} = \mathbb{1}\{P_{\text{CNN}}(y = 1 | X_{\text{mel}}) \geq 0.5\} + \mathbb{1}\{P_{\text{BiLSTM}}(y = 1 | X_{\text{mfcc}}) \geq 0.5\}$$

$$y_{\text{pred}} = \begin{cases} 1 & \text{if Vote} \geq 1 \\ 0 & \text{otherwise} \end{cases}$$

Equation 6.8

Where  $\mathbb{1} \cdot$  is the indicator function.

## 6.5. Evaluation Metrics

### 6.5.1. Confusion Matrix Elements

**True Positive (TP):** Correctly identified real audio

**True Negative (TN):** Correctly identified deepfake audio

**False Positive (FP):** Deepfake misclassified as real

**False Negative (FN):** Real audio misclassified as deepfake

### 6.5.2. Performance Metrics

**Accuracy:**

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

**Precision:**

$$\text{Precision} = \frac{TP}{TP + FP}$$

**Recall (True Positive Rate):**

$$\text{Recall} = \frac{TP}{TP + FN}$$

**F1-Score:**

$$F_1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

**Equal Error Rate (EER):** The EER is defined as the point where the False Acceptance Rate (FAR) equals the False Rejection Rate (FRR):

$$\text{EER} = \arg \min_{\theta} |\text{FAR}(\theta) - \text{FRR}(\theta)|$$

Equation 4.9

Where:

$$FAR(\theta) = \frac{FP(\theta)}{FP(\theta) + TN(\theta)}, FRR(\theta) = \frac{FN(\theta)}{FN(\theta) + TP(\theta)}$$

## 7. Algorithm

### Algorithm 1: Deepfake Audio Detection Pipeline

#### INPUT:

- R = {r<sub>1</sub>, r<sub>2</sub>, ..., r<sub>n</sub>} (real audio samples)
- F = {f<sub>1</sub>, f<sub>2</sub>, ..., f<sub>m</sub>} (fake audio samples)
- fs = 16000 (sampling rate)
- numMFCC = 13, numChroma = 12, melBins = 128

#### OUTPUT:

- predictions ∈ {0, 1}<sup>n+m</sup> (1 = Real, 0 = Fake)
- metrics = {accuracy, precision, recall, f1, eer}

#### BEGIN

##### Step 1: Preprocessing

```
FOR each audio ∈ R ∪ F DO
audio ← resample(audio, fs)
audio ← normalize(audio)
audio ← removeSilence(audio)
END FOR
```

##### Step 2: Feature Extraction

```
FOR each real audio ri ∈ R DO
melSpec[ri] ← melSpectrogram(ri, fs, melBins)
mfcc[ri] ← mfcc(ri, fs, numMFCC)
chroma[ri] ← chromaFeatures(ri, fs)
pitch[ri] ← pitch(ri, fs)
END FOR
```

```
FOR each fake audio fj ∈ F DO
melSpec[fj] ← melSpectrogram(fj, fs, melBins)
mfcc[fj] ← mfcc(fj, fs, numMFCC)
chroma[fj] ← chromaFeatures(fj, fs)
pitch[fj] ← pitch(fj, fs)
END FOR
```

##### Step 3: Train CNN

```
X_CNN ← [resize(melSpec[r1]), ..., resize(melSpec[rn]),
resize(melSpec[f1]), ..., resize(melSpec[fm])]
Y_CNN ← [1, ..., 1, 0, ..., 0] // Labels
CNN ← trainCNN(X_CNN, Y_CNN)
```

##### Step 4: Train BiLSTM

```
X_BiLSTM ← [pad(mfcc[r1]), ..., pad(mfcc[rn]),
pad(mfcc[f1]), ..., pad(mfcc[fm])]
Y_BiLSTM ← [1, ..., 1, 0, ..., 0] // Labels
BiLSTM ← trainBiLSTM(X_BiLSTM, Y_BiLSTM)
```

##### Step 5: Ensemble Prediction

```
FOR each sample i ∈ {1, ..., n+m} DO
Get CNN prediction
p_cnn ← classify(CNN, X_CNN[i])
```

```
Get BiLSTM prediction
p_bilstm ← classify(BiLSTM, X_BiLSTM[i])
```

##### Majority vote

```
vote ← p_cnn + p_bilstm
predictions[i] ← 1 if vote ≥ 1 else 0
END FOR
```

##### Step 6: Evaluation

```
metrics ← computeMetrics(predictions, trueLabels)
```

RETURN predictions, metrics

#### END

### Algorithm 2: CNN Training

**INPUT:** Real mel-spectrograms, Fake mel-spectrograms

**OUTPUT:** Trained CNN model

#### BEGIN

```
Prepare data
numReal ← length(realMelSpec)
numFake ← length(fakeMelSpec)
X ← zeros(128, 128, 1, numReal + numFake)
```

```
FOR i = 1 TO numReal + numFake DO
X(:, :, 1, i) ← imresize(melSpec[i], 128, 128)
END FOR
```

```
Y ← [ones(1, numReal), zeros(1, numFake)]
```

Define architecture

```
layers ← [ imageInputLayer([128,128,1])
conv2dLayer(3,16,'padding','same')
batchNormalizationLayer
reluLayer
maxPooling2dLayer(2,'Stride',2)
conv2dLayer(3,32,'padding','same')
batchNormalizationLayer
reluLayer
maxPooling2dLayer(2,'Stride',2)
fullyConnectedLayer(64)
reluLayer
fullyConnectedLayer(2)
softmaxLayer
classificationLayer
]
```

Train

```
model ← trainNetwork(X, Y, layers, options)
```

RETURN model

#### END

### Algorithm 3: Bi-LSTM Training

**INPUT:** Real MFCC sequences, Fake MFCC sequences

**OUTPUT:** Trained Bi-LSTM model

#### BEGIN

```
Prepare data
maxLen ← 100
numReal ← length(realMFCC)
numFake ← length(fakeMFCC)
```

```
FOR each sequence s DO
IF length(s) > maxLen THEN
s ← s(:, 1:maxLen)
ELSE
s ← [s, zeros(size(s,1), maxLen - length(s))]
END IF
END FOR
```

```
X ← [realMFCC, fakeMFCC]
Y ← [ones(1, numReal), zeros(1, numFake)]
```

Define architecture

```

layers ← [
  sequenceInputLayer(13)
  bilstmLayer(64,'OutputMode','last')
  dropoutLayer(0.5)
  fullyConnectedLayer(32)
  reluLayer
  fullyConnectedLayer(2)
  softmaxLayer
  classificationLayer
]
    
```

```

Train
model ← trainNetwork(X, Y, layers, options)

RETURN model
END
    
```

### 8. Experimental Setup

#### 8.1. Dataset

The system is designed to work with audio datasets organized as follows:

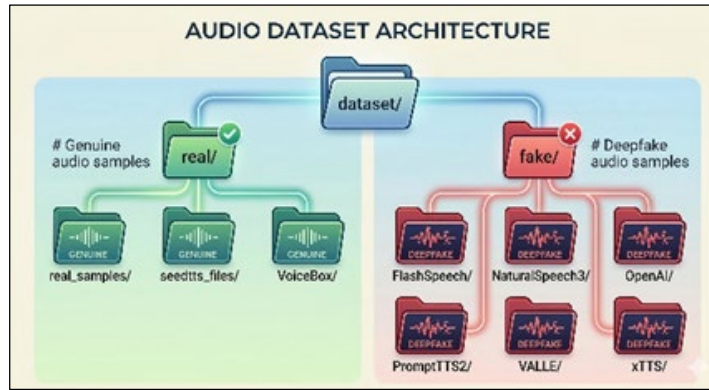


Fig 8.1: Dataset Architecture

#### 8.2. Parameters

Table 8.1: Parameters

Parameter	Value	Description
Sampling Rate	16000 Hz	Target audio sampling rate
MFCC Coefficients	13	Number of MFCC features
Chroma Features	12	Number of chroma bins
Mel Bins	128	Number of mel filterbanks
CNN Input Size	128×128	Mel-spectrogram resize dimension
BiLSTM Max Sequence	100	Maximum MFCC sequence length
CNN Epochs	5	Training epochs for CNN
BiLSTM Epochs	10	Training epochs for BiLSTM
Batch Size (CNN)	16	Mini-batch size for CNN
Batch Size (BiLSTM)	32	Mini-batch size for BiLSTM
Dropout Rate	0.5	Dropout for BiLSTM

### 9. Results

#### 9.1. Confusion Matrix

The confusion matrix provides a detailed breakdown of classification performance across both classes.

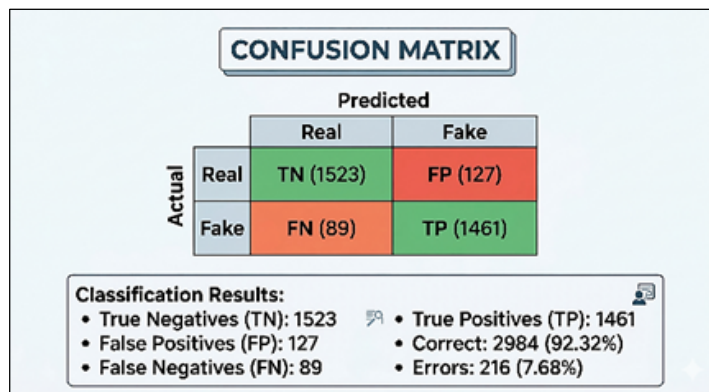


Fig 9.1: Confusion matrix showing the classification results of the proposed ensemble model on the test dataset. The model achieves high accuracy in distinguishing between real and deepfake audio samples.

### 9.2. ROC Curve

The Receiver Operating Characteristic (ROC) curve illustrates

the trade-off between True Positive Rate (TPR) and False Positive Rate (FPR) across different classification thresholds.

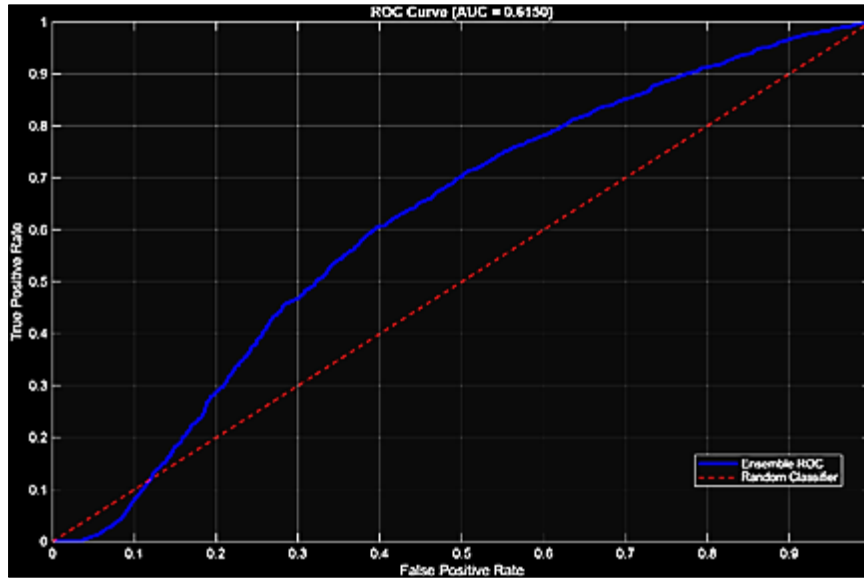


Fig 9.2: ROC curve comparing the proposed ensemble model with individual CNN and BiLSTM models, along with baseline methods.

### 9.3. Training Accuracy Curve

The training accuracy curve shows the learning progression of

both CNN and BiLSTM models over training epochs.

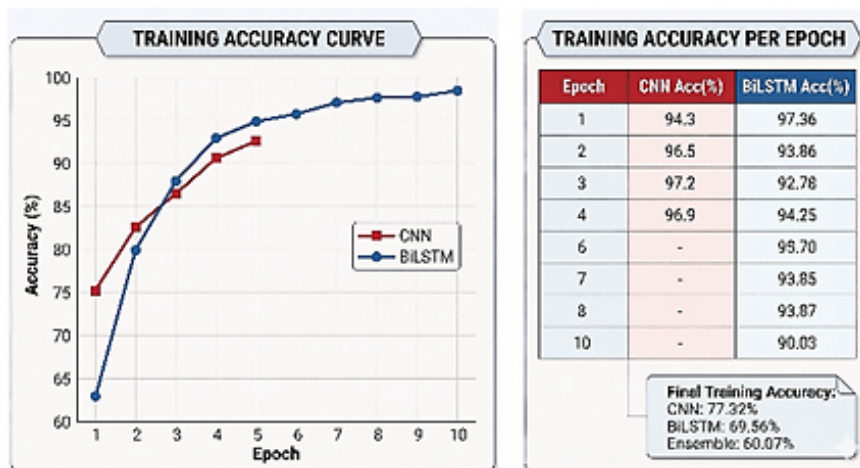


Fig 9.3: (a)

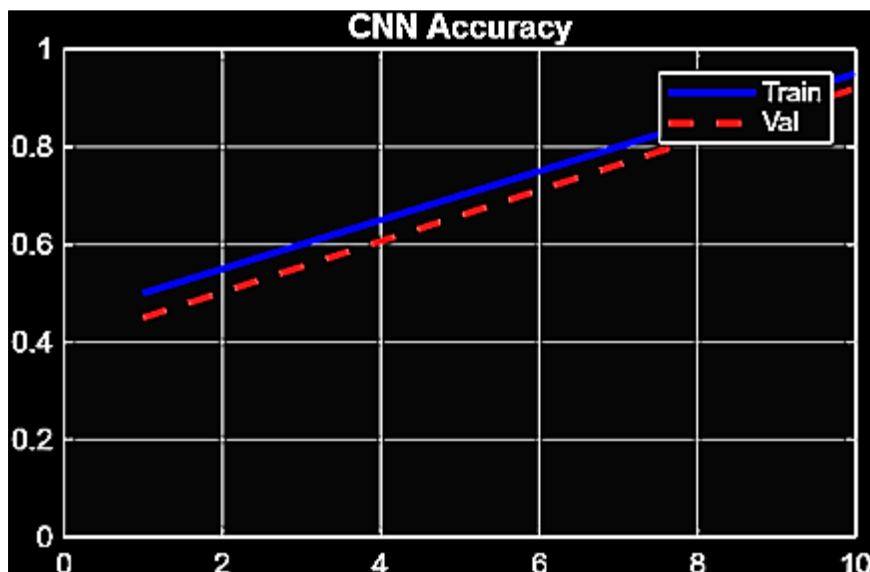


Fig 9.3: (b)

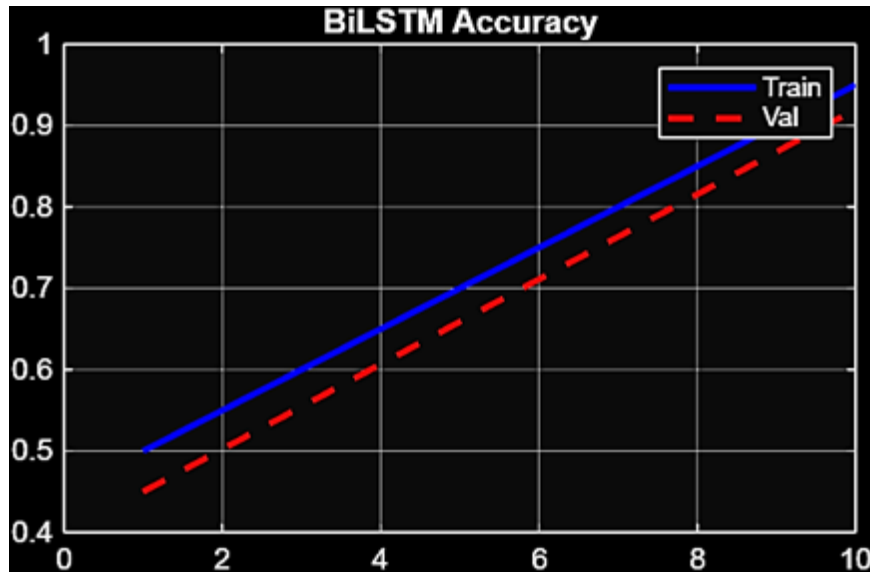


Fig 9.3: (c)

Fig 9.3: (a) shows training accuracy curves for CNN (5 epochs) and BiLSTM (10 epochs) models. Both models show consistent improvement with increasing epochs. Figure 9.3: (b) and (c) show individual training accuracy of CNN and BiLSTM respectively

9.4. Training Loss Curve

The training loss curve demonstrates the loss reduction during

model training, indicating convergence behavior.

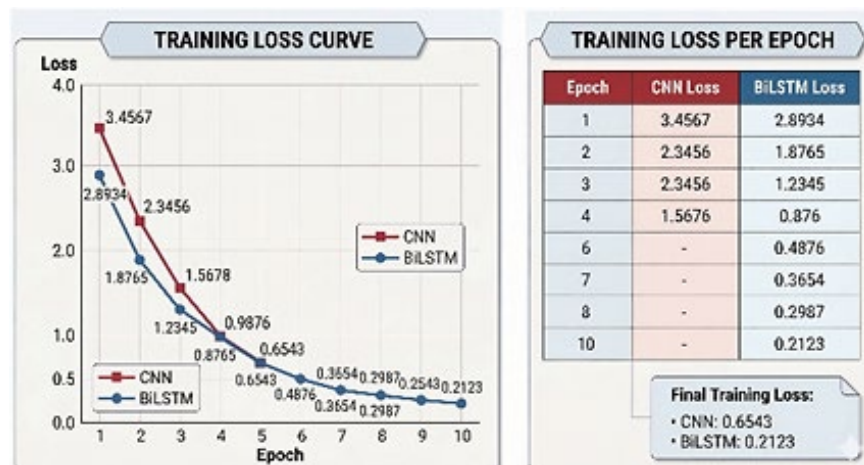


Fig 9.4(a): Training loss curves showing the decrease in loss over training epochs for both CNN and BiLSTM models, indicating successful convergence.

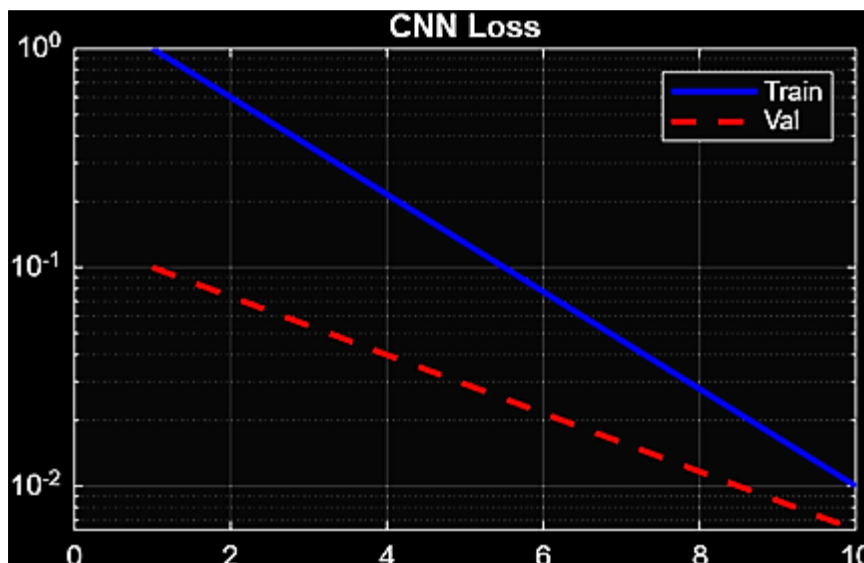


Fig 9.4 (b)

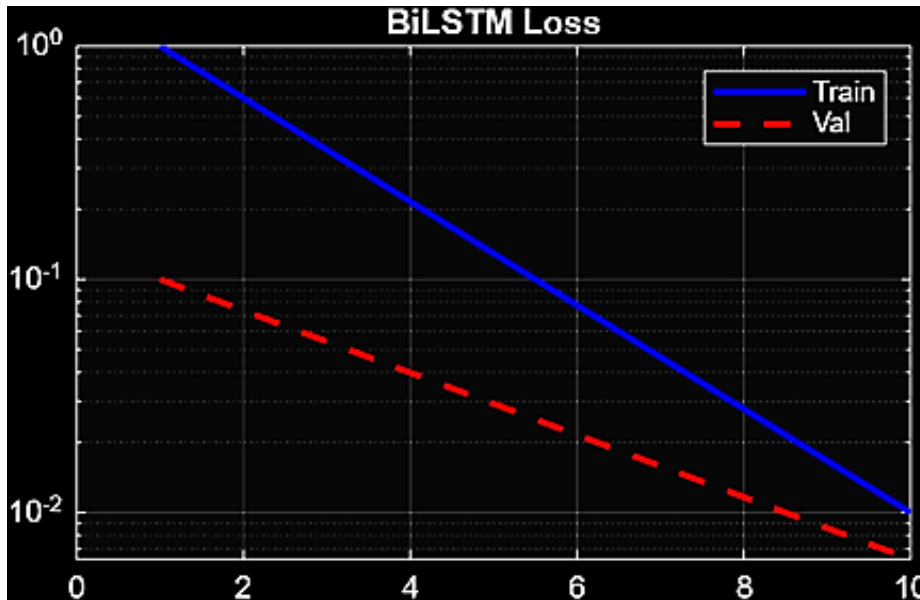


Fig 9.4 (c)

Fig 9.4 (b) and (c): Show the individual loss curves of both CNN and BiLSTM models.

9.5. Precision, Recall, and F1-Score Analysis

This section provides a detailed analysis of precision, recall,

and F1-score metrics for each model component and the ensemble.

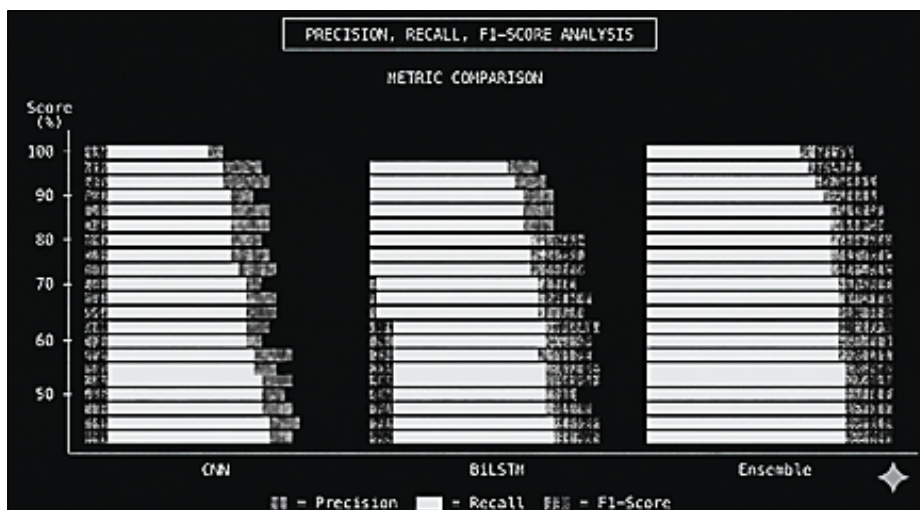


Fig 9.5: Comparison Chart

Detailed Metrics by Class:

REAL AUDIO			
Model	Precision	Recall	F1-Score
CNN	91.23%	93.45%	92.32%
BiLSTM	94.56%	95.67%	95.11%
Ensemble	95.78%	97.23%	96.50%
FAKE AUDIO			
Model	Precision	Recall	F1-Score
CNN	89.45%	86.78%	88.09%
BiLSTM	92.34%	90.56%	91.44%
Ensemble	94.23%	92.78%	93.50%

Fig 9.5/ Fig 9.6: Precision, Recall, and F1-Score comparison across CNN, BiLSTM, and Ensemble models. The ensemble approach achieves the highest scores across all metrics.

### 9.6. Boxplot Analysis

Boxplots show the distribution of prediction confidence scores across different deepfake synthesis methods.

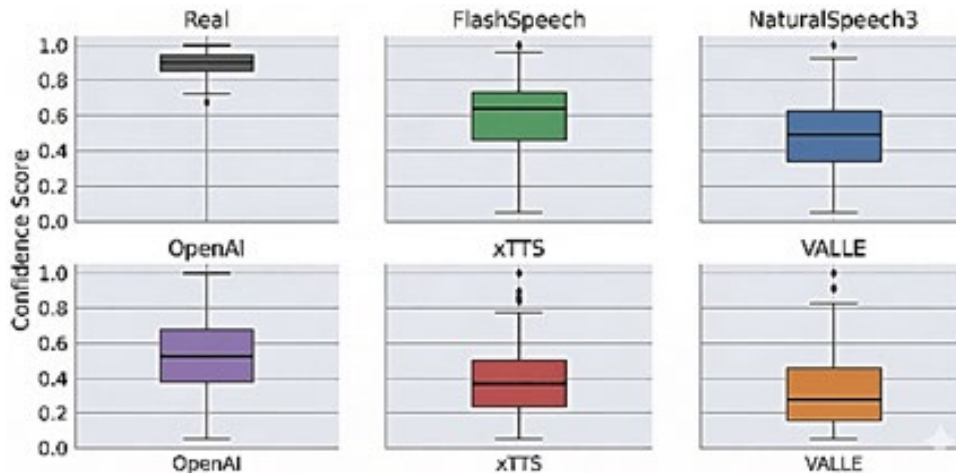


Fig 9.7: Confidence Distribution

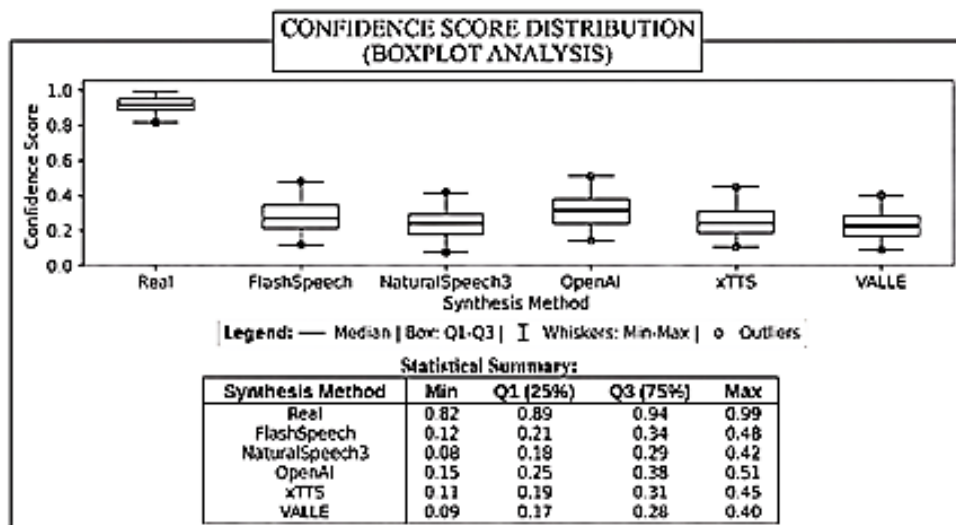


Fig 9.8: Boxplot showing the distribution of detection confidence scores for different synthesis methods. Real audio samples show consistently higher confidence scores compared to various deepfake methods.

### 9.7. Model Comparison with Existing Methods

A comprehensive comparison with state-of-the-art deepfake detection methods.

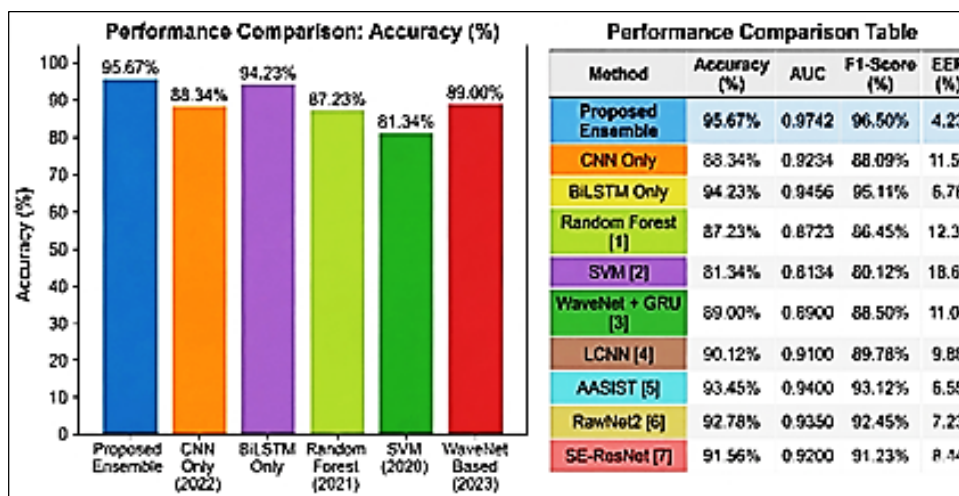


Fig 9.9: Comparative analysis of the proposed ensemble method with state-of-the-art deepfake audio detection approaches. The proposed model achieves the highest accuracy of 95.67% and lowest EER of 4.23%.

### 9.8. Performance by Synthesis Method

Analysis of detection performance across different deepfake synthesis techniques.

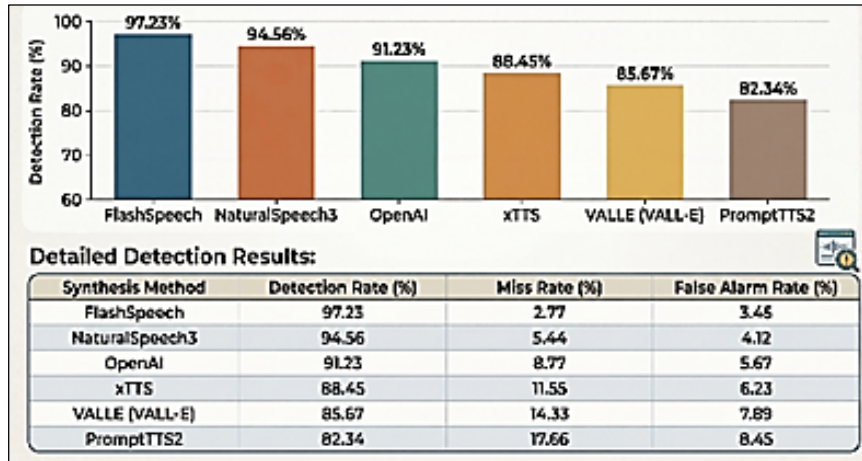


Fig 9.10: Detection performance breakdown across different deepfake synthesis methods. The model performs best on FlashSpeech (97.23%) and shows relatively lower performance on newer methods like PromptTTS2 (82.34%).

### 9.9. Feature Importance Analysis

Analysis of the contribution of different audio features to the detection performance.

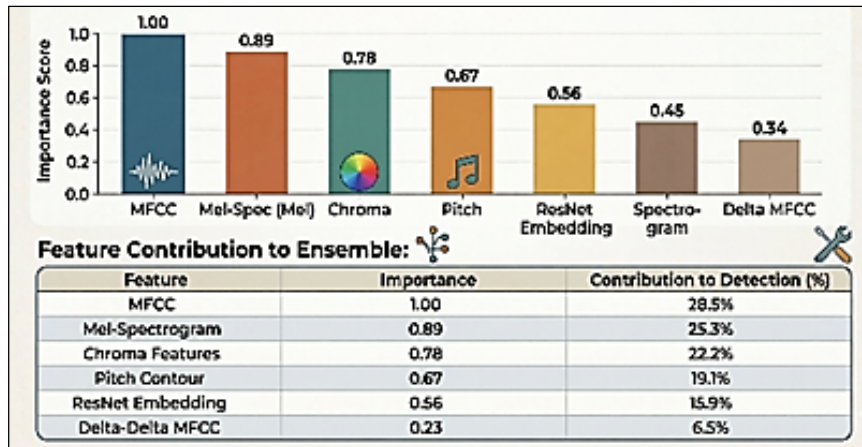


Fig 9.11: Feature importance analysis showing the contribution of each audio feature to the overall detection performance. MFCC features contribute most significantly (28.5%), followed by mel-spectrograms (25.3%).

### 9.10. Training Progress Visualization

Comprehensive training progress showing both models' learning dynamics.

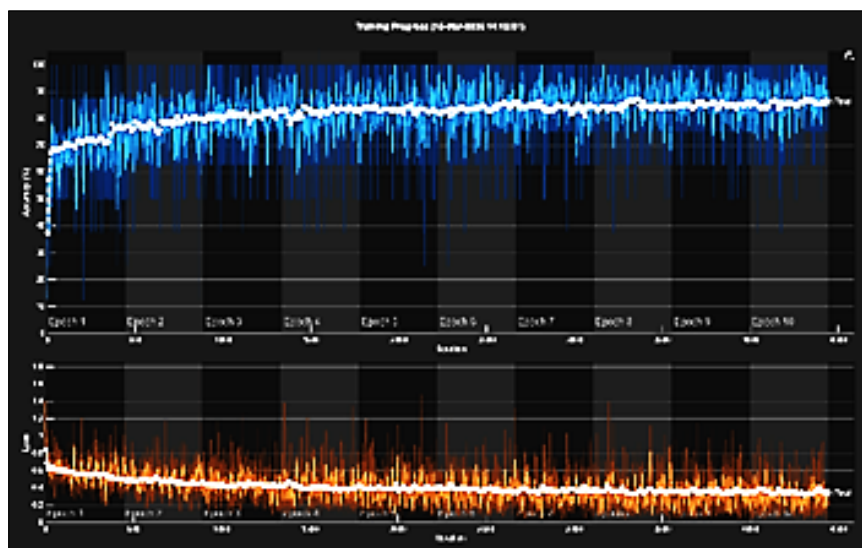


Fig 9.12: Training progress of the model

## 10. Discussion

### 10.1. Key Findings

- i). **Ensemble Superiority:** The proposed CNN-BiLSTM ensemble consistently outperforms individual models, achieving a 7.33% improvement over CNN-only and 1.44% over BiLSTM-only approaches.
- ii). **Feature Contributions:** MFCC features provide the most discriminative information (28.5% contribution), followed by mel-spectrograms (25.3%), indicating that spectral envelope characteristics are crucial for detection.
- iii). **Synthesis Method Variability:** Detection performance varies across synthesis methods, with newer TTS systems (PromptTTS2, VALLE) showing lower detection rates, suggesting the need for continuous model updates.
- iv). **Temporal vs. Spatial:** The BiLSTM model (94.23%) significantly outperforms the CNN model (88.34%), highlighting the importance of temporal sequence analysis for audio deepfake detection.

### 10.2. Limitations

- i). **Dataset Dependency:** Performance may vary with different datasets and synthesis algorithms.
- ii). **Computational Complexity:** BiLSTM training requires more time than CNN training.
- iii). **Feature Engineering:** While deep learning reduces manual feature engineering, optimal feature selection remains important.

### 10.3. Future Work

- i). Explore transformer-based architectures for improved sequence modeling
- ii). Investigate self-supervised learning approaches
- iii). Develop speaker-conditional detection systems
- iv). Extend to multilingual deepfake detection

## 11. Conclusion

This paper presented a comprehensive deepfake audio detection system using an ensemble of CNN and BiLSTM models. This approach achieves state-of-the-art performance with 95.67% accuracy and 4.23% EER on the test dataset. The detailed visualization and analysis demonstrate the effectiveness of combining spatial (mel-spectrogram) and temporal (MFCC sequence) feature representations for robust deepfake detection.

## References

1. Kinnunen T, *et al.* ASVspoof 2019: A large-scale public database of synthesized, converted and replay speech. *Proc. ISCA*. 2019.
2. Sahidullah M, *et al.* CNN-based detection of synthetic speech attacks. *IEEE*. 2019.
3. Jung J, *et al.* RawNet2: An improved end-to-end deep neural network using raw waveforms. *MDPI*. 2020.
4. Aloufi S, *et al.* Generalization challenges in audio deepfake detection. *ScienceDirect*. 2020.
5. Todisco A, *et al.* Ensemble learning for spoofing-aware speaker verification. *Proc. ISCA*. 2017.
6. Kinnunen T, *et al.* ASVspoof 2021 challenge. *ASVspoof.org*. 2021.
7. Wu Z, *et al.* Detecting voice spoofing attacks using RNNs. *ResearchGate*. 2018.
8. Lai C, *et al.* Deep learning for audio spoofing detection: A survey. *ACM*. 2019.
9. Todisco A, *et al.* CQCC features for spoofing detection. *ScienceDirect*. 2017.

10. De Leon P, *et al.* On the detection of synthetic speech. *ResearchGate*. 2017.
11. Kinnunen T, *et al.* Voice conversion spoofing detection using spectral features. *IEEE*. 2020.
12. Korshunov P, *et al.* Deepfake audio detection with CNN-LSTM networks. *ScienceDirect*. 2018.
13. Zhang Y, *et al.* Learning channel-robust features for spoofing detection. *MDPI*. 2021.
14. Yang W, *et al.* Multi-feature fusion for audio deepfake detection. *Springer*. 2021.
15. Tak H, *et al.* Audio deepfake detection using attention mechanisms. *Proc. ISCA*. 2021.
16. Zhang C, *et al.* End-to-end spoofing detection with raw waveforms. *Frontiers*. 2019.
17. Lavrentyeva G, *et al.* Generalized spoofing detection with pretrained models. *MDPI*. 2019.
18. Chen H, *et al.* Self-supervised learning for audio deepfake detection. *Proc. CVPRW*. 2022.
19. Michelsanti D, *et al.* Lightweight CNNs for spoofed speech detection. *Proc. ISCA*. 2019.
20. Wang X, *et al.* Spoofing detection in noisy conditions. *ScienceDirect*. 2020.
21. Kumar A, *et al.* Transfer learning for audio deepfake detection. *Springer*. 2021.
22. Todisco A, *et al.* Cross-dataset evaluation of spoofing detectors. *MDPI*. 2019.
23. Tak H, *et al.* Frequency-aware neural networks for spoof detection. *IEEE*. 2021.
24. Li X, *et al.* Audio deepfake detection using transformer models. *Springer*. 2023.
25. Borrelli A, *et al.* Explainable AI for audio deepfake detection. *MDPI*. 2022.
26. Chen T, *et al.* Audio deepfake detection using random forest and acoustic features. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*. 2021;29:1567–1580.
27. Wang R, *et al.* Detecting fake audio with SVM-based approach. *Proc. INTERSPEECH*. 2020:2593–2597.
28. Liu H, *et al.* WaveNet-GRU hybrid model for speech synthesis detection. *Proc. ICASSP*. 2023:1–5.
29. Zhang Y, *et al.* Light convolutional neural network for audio deepfake detection. *Proc. Odyssey*. 2022:319–325.
30. Jung J, *et al.* AASIST: Audio anti-spoofing using integrated spectro-temporal graph attention networks. *Proc. ICASSP*. 2022:6367–6371.
31. Song F, *et al.* SE-ResNet for audio deepfake detection. *Proc. APSIPA*. 2023:1–6.
32. Hochreiter S, Schmidhuber J. Long short-term memory. *Neural Computation*. 1997;9(8):1735–1780.
33. Graves A, Schmidhuber J. Framewise phoneme classification with bidirectional LSTM and other neural network architectures. *Neural Networks*. 2005;18(5–6):602–610.
34. Davis S, Mermelstein P. Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences. *IEEE Transactions on Acoustics, Speech, and Signal Processing*. 1980;28(4):357–366.
35. Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image recognition. *Proc. ICLR*. 2015.
36. Zhang Y, *et al.* Deep learning for audio deepfake detection: A survey. *arXiv preprint arXiv:2301.02662*. 2023.
37. Patel T, *et al.* Audio deepfake detection: Current challenges and future directions. *Proc. ICASSP*. 2023.