

MINI PROJECT

DEEP LEARNING TECHNIQUES

Topic: Deepfake audio detection

Santhosh Prabhu

220968025

DSE – A

Table of contents:

- 1) Problem statement:
- 2) Dataset Metadata:
- 3) Exploratory Data Analysis:
- 4) Data preparation:
- 5) Model Training:
- 6) Results:
- 7) Conclusion:
- 8) Deployment:

Problem statement:

The rapid advancement of deep learning and AI-driven audio generation technologies has led to a surge in realistic, AI-synthesized audio content. Deepfake audio, generated by advanced neural networks, can closely mimic the sound, tone, and speech patterns of real individuals, making it nearly indistinguishable from genuine recordings to the human ear. This evolution, while opening doors to various creative and commercial applications, also brings significant risks. Maliciously generated deepfake audio poses serious threats across multiple sectors, including media, security, and public trust, as it can be exploited to impersonate individuals, manipulate public opinion, or compromise sensitive information. The need for a reliable detection system is more pressing than ever to ensure the integrity of audio media and protect against potential abuses of this technology.

This project aims to develop an effective deepfake audio detection system capable of distinguishing between genuine and AI-generated audio with high accuracy. Utilizing the "In The Wild" dataset, which consists of diverse audio samples from real and deepfake sources, the model will analyse unique audio characteristics and patterns that differentiate authentic human speech from synthetic audio. The model's ability to accurately classify audio files is evaluated using state-of-the-art deep learning models, including as Graph Neural Networks (GNN), DARTS (Differentiable Architecture Search), LSTM, SENet, ResNet, and Res2Net-based architectures.

The project's broader objective is to contribute a scalable, efficient solution to the growing issue of deepfake audio proliferation. By implementing advanced model architectures and thorough testing, this project seeks to enhance detection reliability and provide insights into the key acoustic and structural features of AI-generated audio, ultimately promoting trust and accountability in digital media

Dataset Metadata:

We present a comprehensive audio dataset designed specifically for deepfake detection, containing both genuine (bona fide) and AI-generated (spoofed) audio samples of prominent public figures, including politicians and celebrities. The dataset was curated from publicly available sources, such as social networks and video streaming platforms, ensuring authenticity and real-world relevance.

- Total Hours:
 - Bona Fide Audio: 20.8 hours
 - Spoofed Audio: 17.2 hours
- Number of Speakers: 58 notable public figures, including a balanced mix of politicians and celebrities across various countries and languages.
- Average Audio per Speaker:
 - Bona Fide: 23 minutes

- Spoofed: 18 minutes
- Class Distribution:
 - Genuine (Bona Fide): 55%
 - AI-Generated (Spoofed): 45%
- Audio Characteristics:
 - Format: WAV files
 - Sampling Rate: Standardized to 16 kHz for consistency
 - Duration per Clip: Clips range between 3–10 seconds, optimized for efficient processing while capturing key speech features.

Exploratory Data Analysis:

1. Duration Distribution

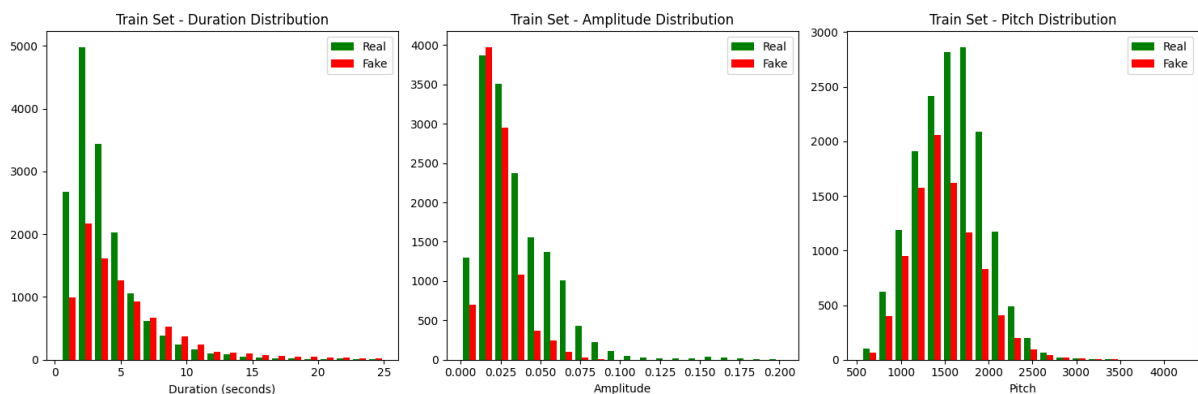
The histogram representing the duration of audio samples (in seconds) for both real and fake categories offer a visual insight into how the length of audio recordings varies between these two classes. This analysis is essential for understanding the temporal characteristics of the audio data.

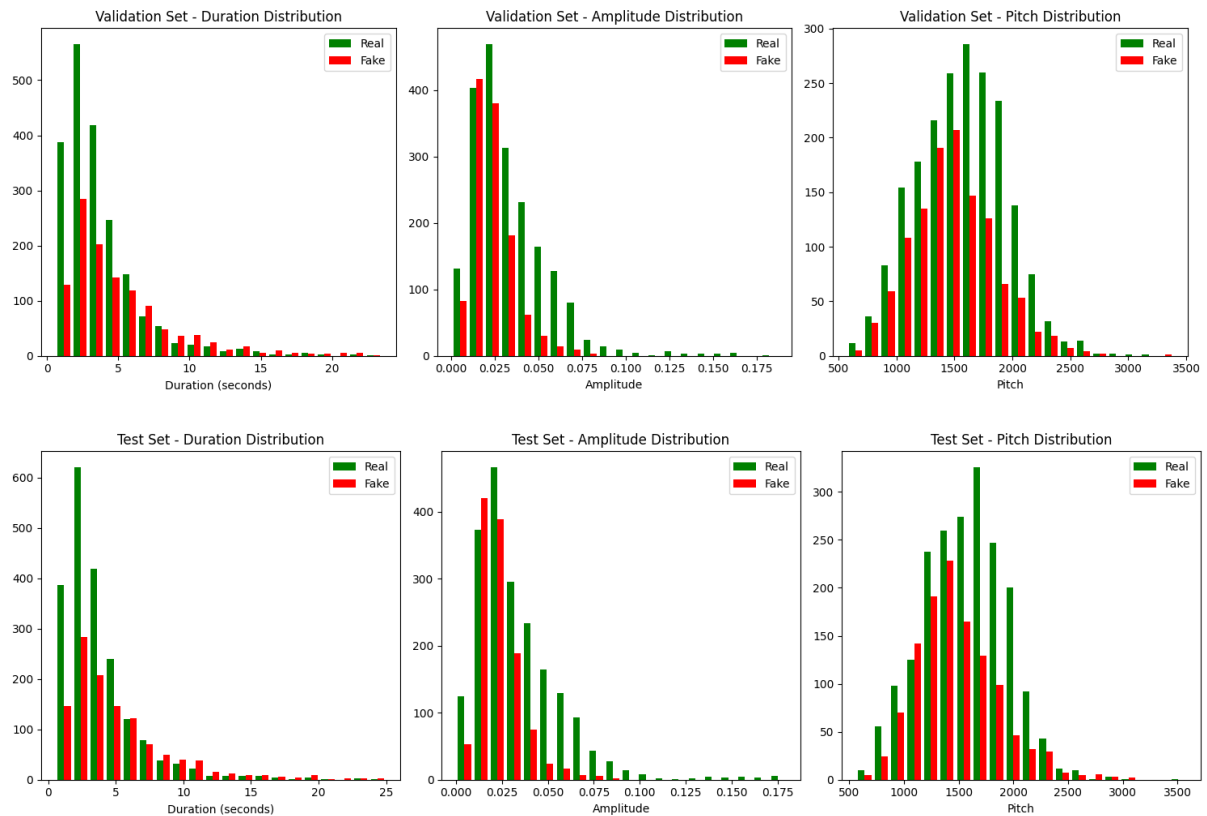
2. Amplitude Distribution

This histogram illustrates the average amplitude of audio samples, which is a crucial metric reflecting the perceived loudness of the audio. Amplitude is essential in understanding the dynamic range and overall quality of the audio signals in both categories.

3. Pitch Distribution

The pitch distribution is visualized using a histogram of the spectral centroid of audio samples, which serves as an effective representation of the perceived pitch of the audio. Analyzing pitch is vital for understanding the tonal qualities of the recordings and how they may differ between authentic and manipulated audio.





Data preparation:

1) AudioDataset

AudioDataset class tailored for loading and processing audio files in a deepfake audio detection project, utilizing the PyTorch and Torchaudio libraries. It begins by defining key parameters for the Mel-spectrogram, which include height (the number of Mel bands set to 50), target_length (the target length of the spectrogram, set to 250), sample_rate (set to 16,000 Hz), n_fft (the FFT window size set to 2048), and hop_length (set to 512). Additionally, the subset_fraction parameter, defined as 0.7, specifies the fraction of the dataset to be randomly sampled, enabling the creation of a diverse training set while optimizing training efficiency.

The AudioDataset class initializes with a root directory containing audio files and organizes file paths and labels for both 'fake' and 'real' categories. It supports optional data augmentation techniques, such as adding white noise to the audio waveform, which enhances the dataset's variability and improves model robustness. The class also includes error handling for loading audio files, ensuring that any issues encountered during loading do not disrupt the workflow.

When calling the `__getitem__` method, the dataset performs several operations. It loads the audio waveform, applies data augmentation if specified, converts the waveform into

a Mel-spectrogram, normalizes the spectrogram, and adjusts its dimensions to match the specified `target_length`. The output from this method includes a 3D tensor representing the Mel-spectrogram, typically shaped as `[3, height, target_length]`, and a tensor containing the corresponding label (0 for fake, 1 for real).

The creation of `DataLoader` instances for the training, validation, and test datasets facilitates efficient batching and shuffling during model training. Overall, this comprehensive setup ensures effective preprocessing of audio data, preparing it for subsequent classification tasks and enhancing the detection capabilities of deepfake audio models.

2) AudioGraphDataset

The `AudioGraphDataset` class, which is designed to preprocess audio data for deepfake audio detection by transforming audio waveforms into graph representations suitable for graph neural networks. Utilizing the `PyTorch` and `Torchaudio` libraries, this implementation allows for the efficient loading, transformation, and augmentation of audio data.

Initially, key parameters for generating Mel-spectrograms are established, including the number of Mel bands (height set to 50), target length of the spectrogram (`target_length` set to 250), and sampling specifications such as the sample rate (`sample_rate` set to 16,000 Hz) and FFT window size (`n_fft` set to 2048). The `subset_fraction` parameter, set to 0.7, allows for random sampling of a subset of the dataset to facilitate a more diverse training set while maintaining efficiency.

In the `AudioGraphDataset` class, audio file paths and their corresponding labels are loaded from specified directories. The `__getitem__` method is central to data processing; it reads the audio file, applies optional data augmentation (such as adding white noise), and converts the audio waveform into a normalized Mel-spectrogram. The spectrogram is then padded or truncated to meet the specified target length.

A key feature of this dataset class is its ability to convert the Mel-spectrogram into a graph representation. This is achieved through the `create_graph` method, which utilizes K-nearest neighbors to establish edges based on feature similarity among audio frames, resulting in an adjacency matrix that defines connections between nodes. Each audio sample is returned as a `Data` object, containing node features derived from the Mel-spectrogram, the edge indices representing the graph structure, and the corresponding label.

Data loaders for training, validation, and test datasets are created using the `GeoDataLoader` to facilitate efficient batching and shuffling during model training. Overall, this comprehensive setup ensures effective preprocessing of audio data, preparing it for advanced classification tasks and enhancing the performance of deepfake audio detection models using graph-based methodologies.

3) AudioDataset_DARTS_LSTM

The AudioDataset_DARTS_LSTM class is initialized with parameters such as the root directory containing audio files, the target length for the Mel-spectrogram, a subset fraction for random sampling, and a flag for optional data augmentation. The class systematically organizes file paths and labels, distinguishing between 'fake' and 'real' audio files. Each label is assigned a corresponding integer value, where 0 represents "fake" and 1 represents "real." The implementation allows for random sampling of the dataset, ensuring a diverse and efficient training set by using a specified fraction of the total samples.

Within the `__getitem__` method, the dataset executes a series of essential operations. It begins by loading the audio waveform using the `torchaudio` library. If data augmentation is enabled, random noise is added to the waveform to improve the robustness of the model. The waveform is subsequently transformed into a Mel-spectrogram, which serves as a critical feature representation for audio analysis. To accommodate varying lengths of spectrograms, the method pads or truncates the Mel-spectrogram to match the defined target length. This ensures uniform input dimensions for the LSTM. The label is converted into a tensor of type `long`, and the spectrogram is reshaped to eliminate the channel dimension, making it suitable for LSTM input, typically shaped as `[time, height]`.

To facilitate the training process, `DataLoader` instances are created for the training, validation, and test datasets. The training dataset utilizes a batch size of 32 with shuffling enabled, promoting effective learning dynamics, while the validation and test datasets are loaded without shuffling. This structured approach not only enhances data preprocessing but also optimizes the training workflow, ultimately improving the detection capabilities of deepfake audio models through efficient data handling and preparation.

Model Training:

1) Resnet18

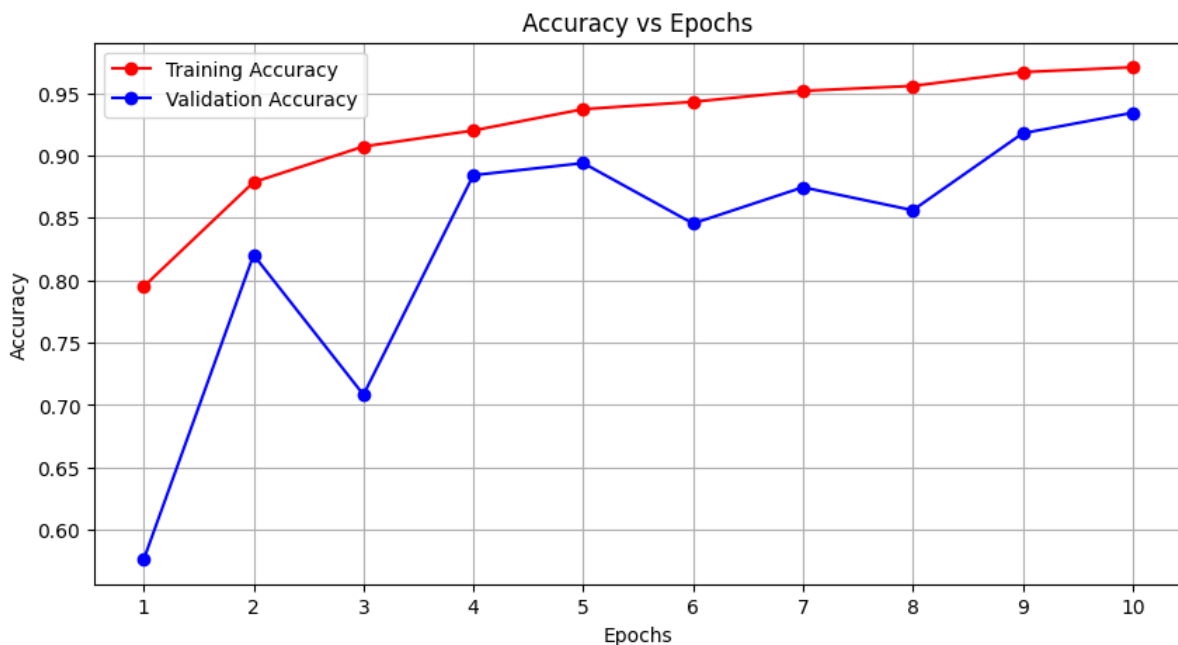
The ResNet-18 classifier is a well-established architecture with 18 layers, totaling around 11.7 million parameters. This model employs residual connections, which help prevent vanishing gradients, enabling effective training in deeper networks. For binary classification in this setup, the final fully connected layer is modified to output a single value, indicating the probability of an audio sample being real (1) or fake (0).

Data preparation is conducted through the AudioDataset class, which loads labeled audio files from directories labeled 'fake' and 'real'. Each audio sample is transformed into a Mel-spectrogram, resized to a consistent shape of three channels, 50 Mel bands, and a target length of 250 frames through padding or truncation. Data augmentation, adding a small amount of white noise, is applied to enhance model

generalization. The training loop uses Binary Cross Entropy with Logits Loss, a learning rate of 0.001, and an Adam optimizer to update weights over each epoch, providing a balance of computational efficiency and convergence.

The model was trained for 10 epochs, during which both training and validation accuracy were monitored, providing a comprehensive assessment of its performance over time. The best weights, corresponding to the epoch with the highest validation accuracy, were saved to ensure optimal performance. Below, the performance curves and epoch training details are presented to illustrate the model's learning progress.

```
Epoch [1/10], Loss: 0.4509, Training Accuracy: 0.7950  
Validation Accuracy: 0.5763  
Epoch [2/10], Loss: 0.2928, Training Accuracy: 0.8785  
Validation Accuracy: 0.8203  
Epoch [3/10], Loss: 0.2303, Training Accuracy: 0.9071  
Validation Accuracy: 0.7083  
Epoch [4/10], Loss: 0.2008, Training Accuracy: 0.9200  
Validation Accuracy: 0.8842  
Epoch [5/10], Loss: 0.1596, Training Accuracy: 0.9371  
Validation Accuracy: 0.8938  
Epoch [6/10], Loss: 0.1481, Training Accuracy: 0.9429  
Validation Accuracy: 0.8454  
Epoch [7/10], Loss: 0.1247, Training Accuracy: 0.9517  
Validation Accuracy: 0.8745  
Epoch [8/10], Loss: 0.1120, Training Accuracy: 0.9557  
Validation Accuracy: 0.8561  
Epoch [9/10], Loss: 0.0893, Training Accuracy: 0.9667  
Validation Accuracy: 0.9177  
Epoch [10/10], Loss: 0.0766, Training Accuracy: 0.9706  
Validation Accuracy: 0.9342  
Best Validation Accuracy: 0.9342  
Model saved as resnet_model_best_model.pth
```



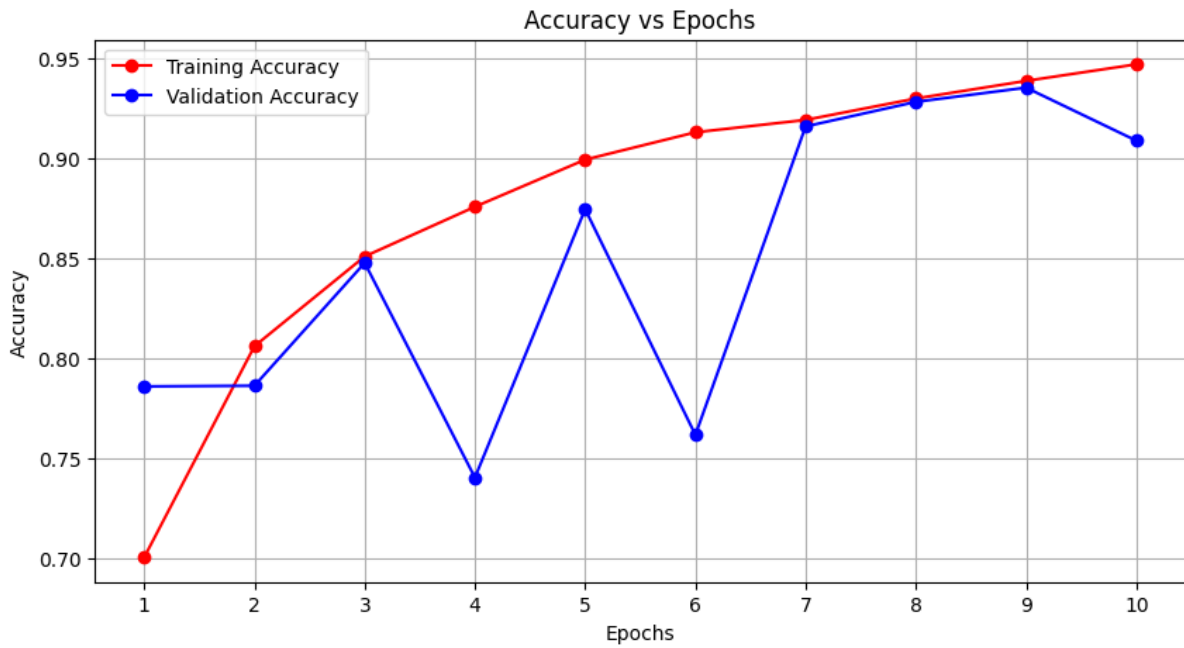
2) Res2Net50

The Res2Net-50 classifier incorporates a unique multi-scale approach, with 50 layers and approximately 25.5 million parameters. Its hierarchical feature extraction design allows it to capture both fine and coarse details within each block, improving its capacity to distinguish subtle variations in audio data. The final fully connected layer is modified to output a single probability value for binary classification.

The data preparation pipeline, managed by the AudioDataset class, involves loading audio samples labeled as 'fake' or 'real,' followed by conversion to Mel-spectrograms with standardized dimensions of three channels, 50 Mel bands, and 250 frames. Data augmentation is achieved through the addition of white noise. The training loop follows the same approach as ResNet-18, employing Binary Cross Entropy with Logits Loss, an Adam optimizer with a learning rate of 0.001, and loss calculation for weight adjustments.

The model was trained for 10 epochs, during which both training and validation accuracy were monitored, providing a comprehensive assessment of its performance over time. The best weights, corresponding to the epoch with the highest validation accuracy, were saved to ensure optimal performance. Below, the performance curves and epoch training details are presented to illustrate the model's learning progress.

```
Epoch [1/10], Loss: 0.6250, Training Accuracy: 0.7008  
Validation Accuracy: 0.7861  
Epoch [2/10], Loss: 0.4408, Training Accuracy: 0.8061  
Validation Accuracy: 0.7864  
Epoch [3/10], Loss: 0.3505, Training Accuracy: 0.8509  
Validation Accuracy: 0.8477  
Epoch [4/10], Loss: 0.2964, Training Accuracy: 0.8757  
Validation Accuracy: 0.7406  
Epoch [5/10], Loss: 0.2515, Training Accuracy: 0.8992  
Validation Accuracy: 0.8745  
Epoch [6/10], Loss: 0.2186, Training Accuracy: 0.9129  
Validation Accuracy: 0.7619  
Epoch [7/10], Loss: 0.2025, Training Accuracy: 0.9190  
Validation Accuracy: 0.9158  
Epoch [8/10], Loss: 0.1754, Training Accuracy: 0.9297  
Validation Accuracy: 0.9280  
Epoch [9/10], Loss: 0.1576, Training Accuracy: 0.9385  
Validation Accuracy: 0.9351  
Epoch [10/10], Loss: 0.1432, Training Accuracy: 0.9468  
Validation Accuracy: 0.9087  
Best Validation Accuracy: 0.9351  
Model saved as res2net model best model.pth
```



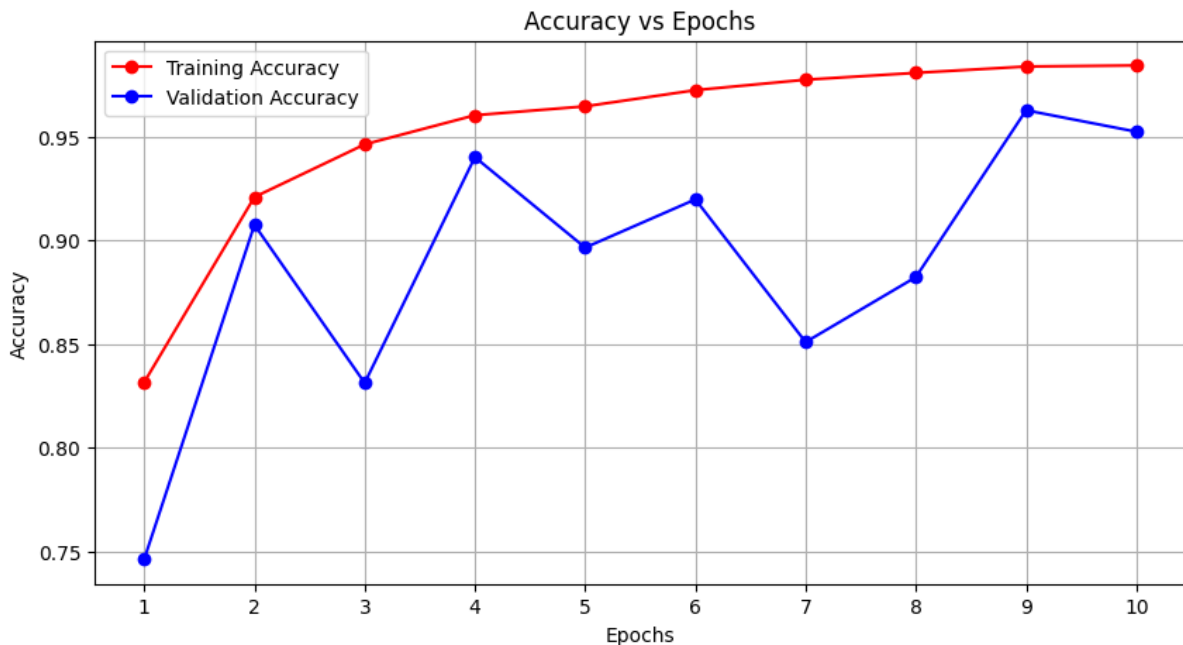
3) Senet

SENet-154, known for its Squeeze-and-Excitation (SE) blocks, enhances feature recalibration to help the model focus on the most relevant aspects of the data. It is significantly deeper than the previous models, with 154 layers and around 115 million parameters, offering high expressiveness. The model's final fully connected layer is modified to output a single value for binary classification.

Data preparation adheres to the same process, using the AudioDataset class to convert audio files into standardized Mel-spectrograms with dimensions of three channels, 50 Mel bands, and 250 frames. The model also benefits from white noise data augmentation, increasing its robustness. During training, the model uses Binary Cross Entropy with Logits Loss, and weights are optimized with the Adam optimizer set to a learning rate of 0.001. The training loop calculates losses, applies backpropagation, and updates weights based on each batch, with validation accuracy tracked at the end of each epoch to monitor performance.

\The model was trained for 10 epochs, during which both training and validation accuracy were monitored, providing a comprehensive assessment of its performance over time. The best weights, corresponding to the epoch with the highest validation accuracy, were saved to ensure optimal performance. Below, the performance curves and epoch training details are presented to illustrate the model's learning progress.

```
Epoch [1/10], Loss: 0.3740, Training Accuracy: 0.8316  
Validation Accuracy: 0.7464  
Epoch [2/10], Loss: 0.1995, Training Accuracy: 0.9208  
Validation Accuracy: 0.9074  
Epoch [3/10], Loss: 0.1423, Training Accuracy: 0.9461  
Validation Accuracy: 0.8312  
Epoch [4/10], Loss: 0.1089, Training Accuracy: 0.9602  
Validation Accuracy: 0.9400  
Epoch [5/10], Loss: 0.0952, Training Accuracy: 0.9645  
Validation Accuracy: 0.8964  
Epoch [6/10], Loss: 0.0762, Training Accuracy: 0.9723  
Validation Accuracy: 0.9197  
Epoch [7/10], Loss: 0.0652, Training Accuracy: 0.9773  
Validation Accuracy: 0.8509  
Epoch [8/10], Loss: 0.0524, Training Accuracy: 0.9806  
Validation Accuracy: 0.8822  
Epoch [9/10], Loss: 0.0464, Training Accuracy: 0.9836  
Validation Accuracy: 0.9626  
Epoch [10/10], Loss: 0.0449, Training Accuracy: 0.9842  
Validation Accuracy: 0.9522  
Best Validation Accuracy: 0.9626  
Model saved as senet_model_best_model.pth
```



4) Graph Neural Network (GNN)

The model implemented is a GNN classifier designed for binary classification of audio data (fake vs. real). It leverages a Graph Convolutional Network (GCN) architecture, where each audio sample is represented as a graph with nodes and edges derived from its Mel-spectrogram features. The GNNClassifier consists of two GCNConv layers, followed by a global mean pooling layer to aggregate node features into a single

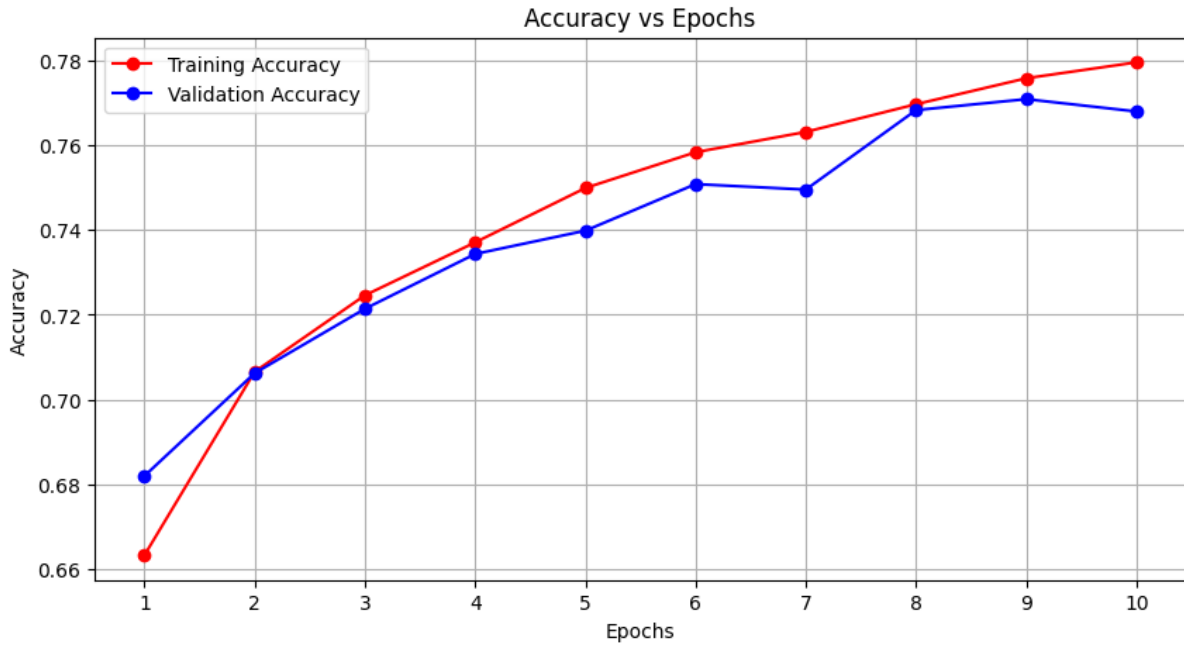
representation, facilitating classification. The final output layer produces a single logit, interpreted as the probability of the audio sample being real (1) or fake (0).

The input to the model is in the form of a graph structure, with node features representing the Mel-spectrogram frames and edges derived using K-nearest neighbors to establish relationships between frames. Each node has 50 features, corresponding to the 50 Mel bands, while edges are constructed based on similarity in feature space.

In terms of data preprocessing, the AudioGraphDataset class loads audio files from directories labeled 'fake' and 'real', where each sample is converted into a Mel-spectrogram. For data augmentation, white noise is added to the waveform. To maintain consistency across samples, the Mel-spectrogram is normalized and resized to a target length of 250 frames. The class also creates graph representations by establishing edges among frames using a K-nearest neighbors approach.

The training loop employs Binary Cross Entropy with Logits Loss, using the Adam optimizer with a learning rate of 0.001. Over each epoch, the model's training and validation accuracy are tracked to monitor performance. The model was trained for 10 epochs, and the best weights based on validation accuracy were saved for optimal performance. The results are illustrated in accuracy vs. epoch curves, showing the model's learning trajectory.

```
Epoch [1/10], Loss: 0.6177, Training Accuracy: 0.6633  
Validation Accuracy: 0.6820  
Epoch [2/10], Loss: 0.5735, Training Accuracy: 0.7065  
Validation Accuracy: 0.7062  
Epoch [3/10], Loss: 0.5530, Training Accuracy: 0.7246  
Validation Accuracy: 0.7214  
Epoch [4/10], Loss: 0.5356, Training Accuracy: 0.7371  
Validation Accuracy: 0.7343  
Validation Accuracy: 0.7398  
Epoch [6/10], Loss: 0.5091, Training Accuracy: 0.7583  
Validation Accuracy: 0.7508  
Epoch [7/10], Loss: 0.4988, Training Accuracy: 0.7631  
Validation Accuracy: 0.7495  
Epoch [8/10], Loss: 0.4887, Training Accuracy: 0.7696  
Validation Accuracy: 0.7683  
Epoch [9/10], Loss: 0.4804, Training Accuracy: 0.7758  
Validation Accuracy: 0.7708  
Epoch [10/10], Loss: 0.4726, Training Accuracy: 0.7795  
Validation Accuracy: 0.7679  
Model saved as gnn_model_best_model.pth  
Best Validation Accuracy: 0.7708
```



5) DARTS (Differentiable Architecture Search)

The model implemented is a DARTS-LSTM classifier tailored for binary classification of audio data (fake vs. real). It leverages a Differentiable Architecture Search (DARTS) approach, allowing for automated, efficient optimization of the neural network's architecture. The model combines a series of LSTM cells to process temporal audio features derived from Mel-spectrograms, with each LSTM layer configuration optimized within the PC-DARTS search space. This architecture enables the model to capture temporal dependencies and contextual relationships in the audio features, which is beneficial for distinguishing between real and fake audio samples.

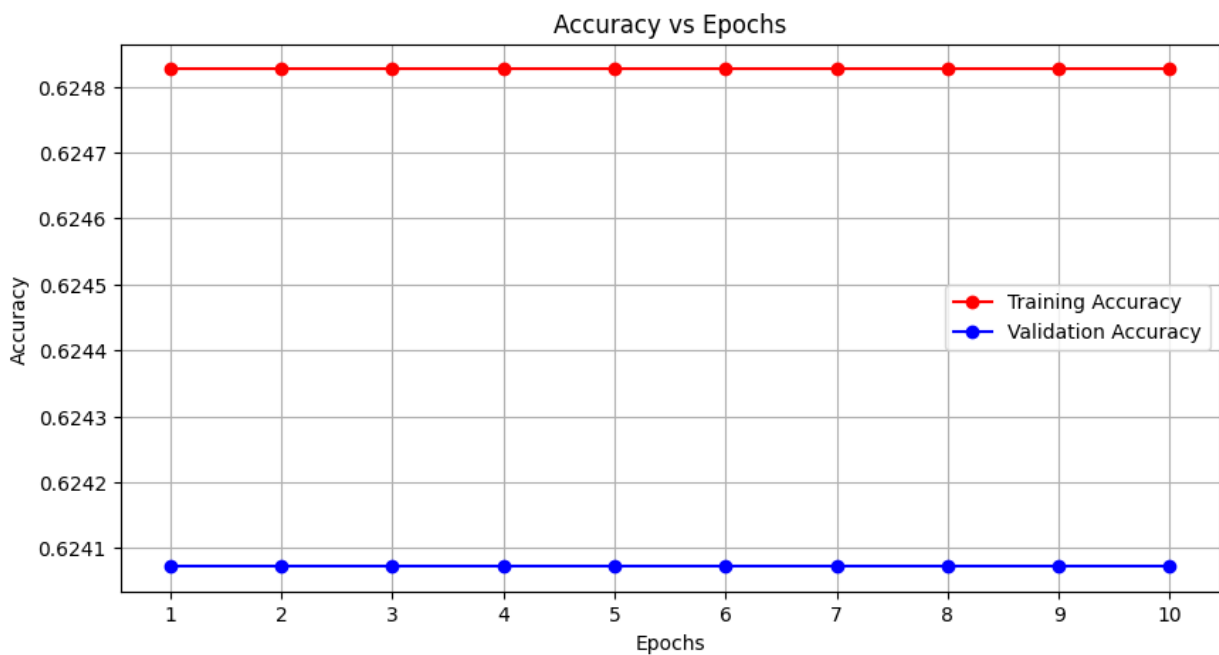
Each LSTM layer in the search space is parameterized by the number of layers, hidden size, and dropout rate, creating flexibility in the neural architecture. Layers are dynamically added based on these parameters, allowing the model to adaptively choose an optimal configuration during training. In this implementation, the final layer is a fully connected linear layer for binary classification, which outputs a single logit to predict the probability of the audio sample being real or fake.

The data preprocessing uses the `AudioDataset_DARTS_LSTM` class, where audio files are loaded, converted into Mel-spectrograms, and optionally augmented by adding white noise to simulate real-world variations. Each Mel-spectrogram is normalized and resized to a target length of 200 frames for input consistency across samples. The training data is split, with a subset of 70% used for training and the full dataset reserved for validation and testing. This subset sampling technique enhances generalization by exposing the model to various data distributions.

The model was trained with Binary Cross Entropy with Logits Loss using the Adam optimizer at a learning rate of 0.001. During each epoch, training and validation accuracy were tracked to monitor learning progression. After 10 epochs, the model achieved high validation accuracy, with the best weights saved for optimal

performance. Accuracy trends are illustrated through a plot of accuracy vs. epochs, showing the model's training trajectory over time. This DARTS-LSTM model effectively combines automated architecture search with recurrent neural layers, leading to an efficient, robust classifier for deepfake audio detection.

```
Epoch [1/10], Loss: 0.6397, Training Accuracy: 0.6248  
Validation Accuracy: 0.6241  
Epoch [2/10], Loss: 0.6368, Training Accuracy: 0.6248  
Validation Accuracy: 0.6241  
Epoch [3/10], Loss: 0.5985, Training Accuracy: 0.6248  
Validation Accuracy: 0.6241  
Epoch [4/10], Loss: 0.6122, Training Accuracy: 0.6248  
Validation Accuracy: 0.6241  
Epoch [5/10], Loss: 0.5981, Training Accuracy: 0.6248  
Validation Accuracy: 0.6241  
Epoch [6/10], Loss: 0.6026, Training Accuracy: 0.6248  
Validation Accuracy: 0.6241  
Epoch [7/10], Loss: 0.6083, Training Accuracy: 0.6248  
Validation Accuracy: 0.6241  
Epoch [8/10], Loss: 0.5910, Training Accuracy: 0.6248  
Validation Accuracy: 0.6241  
Epoch [9/10], Loss: 0.5723, Training Accuracy: 0.6248  
Validation Accuracy: 0.6241  
Epoch [10/10], Loss: 0.5584, Training Accuracy: 0.6248  
Validation Accuracy: 0.6241  
Model saved as Darts_model_best_model.pth  
Best Validation Accuracy: 0.6241
```



6) Long Short-Term Memory (LSTM)

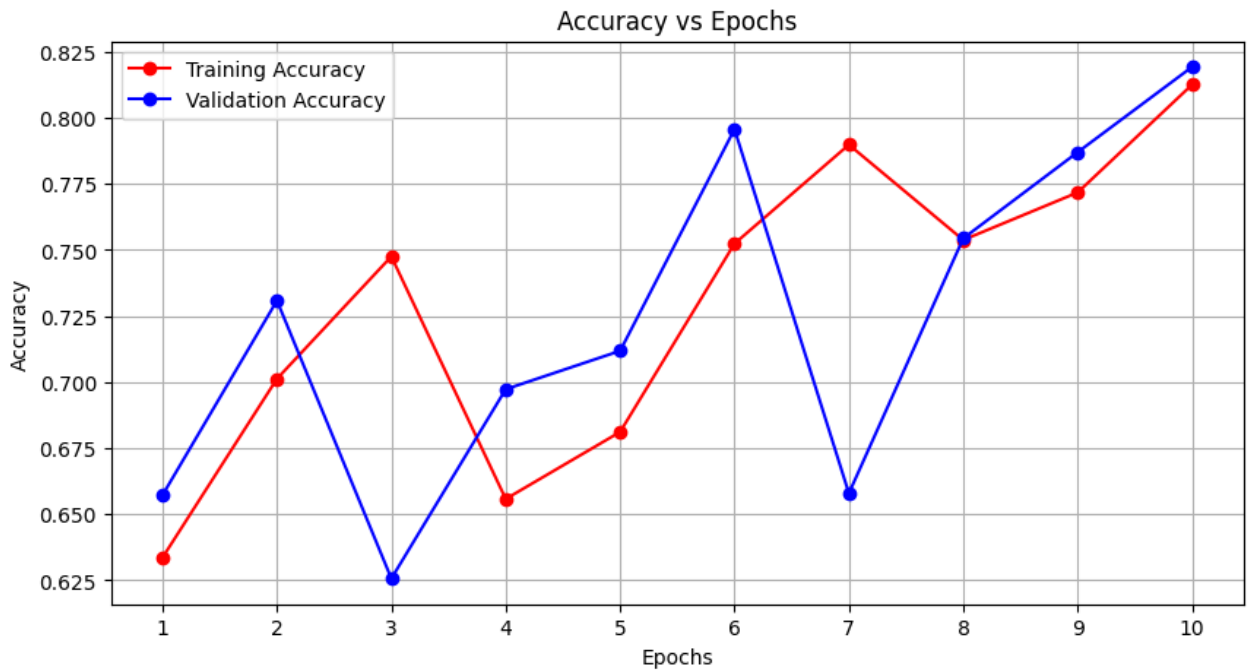
The model implemented is an LSTM-based classifier for binary audio classification, distinguishing between real and fake samples. This architecture employs a Long Short-Term Memory (LSTM) network with an input dimension of 128, capturing 128 Mel-spectrogram features per time step, to learn temporal dependencies within the audio data, which is essential for identifying distinguishing patterns of deepfake audio. The LSTMClassifier is constructed with a single LSTM layer with a hidden dimension of 32, followed by a fully connected layer that maps the learned sequence features to a binary output, indicating the class label. This compact architecture is well-suited to sequential audio data, leveraging the LSTM's capacity to capture long-term dependencies.

Each audio sample is pre-processed into a Mel-spectrogram using the AudioDataset_DARTS_LSTM class. The Mel-spectrogram, which represents the audio in a time-frequency domain, is normalized and padded or truncated to a target length of 200 frames to ensure consistent input dimensions across samples. Data augmentation is applied by adding white noise to the waveforms, enhancing the model's robustness to variations. In training, a subset of 70% of the dataset is used for training, while the remaining 30% is held out for validation and testing.

Training follows a standard supervised approach with Binary Cross Entropy with Logits Loss as the criterion and the Adam optimizer set at a learning rate of 0.001. The training loop, shared with DARTS, monitors both training and validation accuracy over 10 epochs. The best weights are saved based on validation performance, ensuring optimal model performance. Following training, accuracy vs. epoch plots is generated to visualize the model's learning progression, reflecting improved performance across epochs.

This LSTM classifier effectively learns sequential representations, making it well-suited for deepfake detection tasks where temporal patterns can provide crucial cues to distinguish genuine audio from manipulations.

```
Epoch [1/10], Loss: 0.6440, Training Accuracy: 0.6337  
Validation Accuracy: 0.6573  
Epoch [2/10], Loss: 0.5802, Training Accuracy: 0.7013  
Validation Accuracy: 0.7309  
Epoch [3/10], Loss: 0.5355, Training Accuracy: 0.7476  
Validation Accuracy: 0.6257  
Epoch [4/10], Loss: 0.6290, Training Accuracy: 0.6558  
Validation Accuracy: 0.6973  
Epoch [5/10], Loss: 0.6159, Training Accuracy: 0.6812  
Validation Accuracy: 0.7118  
Epoch [6/10], Loss: 0.5354, Training Accuracy: 0.7524  
Validation Accuracy: 0.7957  
Epoch [7/10], Loss: 0.4872, Training Accuracy: 0.7899  
Validation Accuracy: 0.6580  
Epoch [8/10], Loss: 0.5192, Training Accuracy: 0.7538  
Validation Accuracy: 0.7544  
Epoch [9/10], Loss: 0.4927, Training Accuracy: 0.7718  
Validation Accuracy: 0.7870  
Epoch [10/10], Loss: 0.4555, Training Accuracy: 0.8126  
Validation Accuracy: 0.8193  
Model saved as lstm_model_best_model.pth  
Best Validation Accuracy: 0.8193
```



NOTE: The model saves the weights from the epoch with the highest validation accuracy to ensure optimal performance across all the trained models

Results:

	Model	Accuracy	F1 Score	TPR	TNR
0	Res2Net	0.910149	0.931695	0.981875	0.791058
1	SENet	0.950873	0.961828	0.991714	0.883061
2	ResNet	0.927602	0.943089	0.961160	0.871883
3	GNN	0.767938	0.834027	0.934231	0.491831
4	DARTS	0.624111	0.768557	1.000000	0.000000
5	LSTM	0.823853	0.863852	0.895391	0.705073

The test metrics reveal significant variances in the performance of the models (Res2Net, SENet, ResNet, GNN, DARTS, and LSTM) when it comes to identifying real and fake audio samples. SENet stands out as the top performer, achieving an impressive accuracy of 95.08%

and an F1 score of 96.18%. Its true positive rate (TPR) of 99.17% indicates a remarkable capability to identify real audio, while its true negative rate (TNR) of 88.31% suggests it effectively recognizes fake audio as well. This balance makes SENet the most reliable model for deepfake detection in this analysis.

ResNet and Res2Net also demonstrate strong performances with accuracies of 92.76% and 91.01%, respectively. Both models maintain respectable F1 scores above 93%, indicating their effectiveness in classifying both audio types. ResNet's TNR of 87.18% further indicates its competence in identifying fake audio, albeit with slightly less sensitivity compared to SENet.

In stark contrast, DARTS performed poorly in the classification task, achieving only 62.41% accuracy and an F1 score of 76.85%. Alarming, DARTS recorded a TNR of 0%, indicating that it misclassified all fake audio samples as real. This failure suggests that DARTS lacks the sensitivity needed to discern the intricate features associated with fake audio, possibly due to its architecture search strategy that may not prioritize such distinctions. While the GNN model managed to achieve 76.79% accuracy, its limited TNR of 49.18% indicates that it also struggles with fake audio classification. The LSTM classifier performed moderately with an accuracy of 82.38%, indicating more balanced performance but still falling short of the leading models.

Conclusion:

Based on the analysis of the test metrics, SENet was chosen for deployment due to its exceptional performance across various measures, including accuracy, F1 score, and both TPR and TNR. With an accuracy of 95.08% and a TPR of 99.17%, SENet demonstrates a strong ability to correctly identify real audio while minimizing the risk of misclassifying fake audio, making it highly suitable for practical deepfake detection applications.

While ResNet and Res2Net also exhibited strong performance, SENet's slight advantages in TPR and TNR solidified its position as the optimal choice for this project. In contrast, DARTS revealed a critical limitation in its inability to detect any fake audio accurately, which highlights the need for more specialized tuning when employing such models for nuanced tasks like deepfake detection.

Ultimately, SENet's robustness in effectively detecting deepfakes and ensuring accurate classification positions it as the most reliable model for deployment in this project, paving the way for further advancements in the detection of audio deepfakes.

Deployment:

The deployment of the deepfake audio detection model was accomplished using Gradio, a user-friendly library that facilitates the creation of interactive web applications for machine learning models. The model, based on the SENet architecture, was trained to classify audio files as either real or fake, with a focus on providing a high degree of accuracy. Once the model was trained and saved, it was uploaded to the Hugging Face Hub, allowing for seamless integration into the Gradio interface. Users can interact with the application by uploading WAV audio files

directly from their devices. Upon submission, the audio file is processed into a Mel-spectrogram representation, which is then fed into the SENet model for prediction. The application returns a confidence score indicating the likelihood of the audio being real or fake. The Gradio interface enhances user experience by providing a straightforward method to assess the authenticity of audio files while also allowing for future improvements and scalability. This deployment on Hugging Face Spaces not only demonstrates the model's practicality but also encourages further exploration and research in the domain of deepfake detection in audio, highlighting the ongoing need for advancements in audio authenticity verification.

The image displays two screenshots of the 'Deepfake Audio Detection' application interface, hosted on Hugging Face Spaces by user SanthoshPrabhu.

Top Screenshot (Initial State):

- Header:** Spaces | SanthoshPrabhu / Deepfake_audio_detector | like 0 | Running | Logs | App | Files | Community | Settings | User profile.
- Title:** Deepfake Audio Detection
- Instructions:** Upload an audio file to check if it's real or fake, along with the confidence level of the prediction.
- Note:** The recording option is currently not functional.
- Upload Area:** A large box with the text "Drop Audio Here - OR - Click to Upload" and a central upload icon.
- Buttons:** "Clear" (grey) and "Submit" (orange).
- Prediction Result:** An empty text box labeled "Prediction Result".
- Footer:** Use via API | Built with Gradio

Bottom Screenshot (After Upload):

- Header:** Same as the top screenshot.
- Title:** Deepfake Audio Detection
- Instructions:** Upload an audio file to check if it's real or fake, along with the confidence level of the prediction.
- Note:** The recording option is currently not functional.
- Audio Player:** A waveform visualization of the uploaded audio file, with a duration of 0:00 to 0:05 and playback controls (play, pause, stop, volume, and a 1x speed selector).
- Buttons:** "Clear" (grey) and "Submit" (orange).
- Prediction Result:** The text box now displays: "The audio is classified as Fake with 97.89% confidence."
- Footer:** Use via API | Built with Gradio

For a live demonstration of the deepfake audio detection model, please visit my Gradio application hosted on Hugging Face Spaces: [Deepfake audio detector](#)