# Among us

- Detecting spoofed audio samples in a multi-party conversation -

Project Report

Vigneshwar Anandhamurugan

Aalborg University
Cyber Security

# AALBORG UNIVERSITY

## STUDENT REPORT

**Title:**
Among us - Detecting spoofed audio samples in a multi-party conversation

**Theme:**
Scientific Theme

**Project Period:**
Fall Semester 2025

**Project Group:**
1017

**Participant(s):**
Vigneshwar Anandhamurugan

**Supervisor(s):**
Ashutosh Dhar Dwivedi
Hannes Künstner

**Page Numbers:** 42

**Date of Completion:**
June 3, 2025

**Abstract:**

As voice technologies become more prevalent, it becomes vital to improve spoofing detection algorithms especially when multiple speakers are involved. This thesis presents a novel pipeline combining speech separation with spoof detection to address the challenge of identifying spoofed audio samples in overlapped speech. The system uses a SOTA spoof detection model which achieves an EER of 2% on ASVspoof 2019 dataset. However, the EER drops to 24% when tested on a custom- generated mixed audio dataset, as a result of modifications of the audio artifacts by the speech separation process. The system when trained on the custom dataset saw an increase in EER to 22% highlighting the importance of domain adaptation for spoof detection in complex acoustic environments and provides a foundation for future research in real-world scenarios.

# Contents

# Preface

Aalborg University, 2025

---

Vigneshwar Anandhamurugan
vanand23@student.aau.dk

# Abbreviations

A list of the abbreviations used in this report, sorted in alphabetical order:

| Abbreviation | Definition |
|---|---|
| AI | Artificial Intelligence |
| AUC | Area Under the Curve |
| AUROC | Area Under the Receiver Operating Characteristic curve |
| BB | Bonafide-Bonafide |
| CNN | Convolutional Neural Network |
| CQCC | Constant Q Cepstral Coefficients |
| CPU | Central Processing Unit |
| CSV | Comma-Separated Values |
| CUDA | Compute Unified Device Architecture |
| DNN | Deep Neural Network |
| EER | Equal Error Rate |
| FBI | Federal Bureau of Investigation |
| GAN | Generative Adversarial Network |
| GDPR | General Data Protection Regulation |
| GPU | Graphics Processing Unit |
| kHz | Kilohertz |
| LCNN | Light Convolutional Neural Network |
| LR | Learning Rate |
| MFCC | Mel Frequency Cepstral Coefficients |
| ML | Machine Learning |
| RNN | Recurrent Neural Network |
| SB | Spoof-Bonfide |
| SOTA | State Of The Art |
| SS | Spoof-Spoof |
| STFT | Short-Time Fourier Transform |
| US | United States |
| VSDC | Voice Spoofing Detection Corpus (audio dataset) |
| ZCR | Zero-Crossing Rate |

# Chapter 1

# Introduction

Audio Deepfakes are AI generated or edited/synthesized to create fake audio that seems real [1]. They often use sophisticated machine learning algorithms such as deep neural networks (DNNs), generative adversarial networks (GANs), convolutional neural networks (CNNs), and recurrent neural networks (RNNs)[2]. Despite the different architectures of these machine learning algorithms, they primarily process audio clips to discern unique patterns in a person's voice such as vocabulary, accent, and speech patterns, which are then used to create audio clips that resemble the original voice note.

Audio deepfakes have become a growing concern with the recent boom in AI and AI associated deepfake generating techniques. In the United States, FBI has issued warnings about threat actors using sophisticated AI techniques to mimic the voices of senior US officials. These scammers use artificially crafted voice notes to deceive individuals, into providing personally identifiable information, credentials, and/or access to personal accounts. The impersonations appear highly realistic, making it difficult to discern the fraud [3]. On the other hand, an Argentinian woman was deceived by scammers impersonating George Clooney. The victim was convinced that she was in contact with the "real" actor for over six weeks, eventually transferred approximately £10,000 to the fraudsters. [4].

Despite promising strides in the development of techniques that identify spoofed audio samples, most of the work in this realm has been conducted with static audio clips. Such datasets are compromised of audio clips that last anywhere between 2-10 seconds, with little or no attention given to variations in acoustic background and recording conditions. Furthermore, these clips usually belong to a single person. [5] tests the reproducibility of such models in a real-time scenario, such as a Teams group call, and reported that these models failed to accurately classify intermittently induced spoofed audio streams. This report addresses the issue by identifying the possible cause of this drop in performance and proposes possible solutions to address this research gap. This report is organized in the following manner: Background, Related work, Methodology, Results, Discussion and

Conclusion.

## 1.1 Problem Statement

With the increasing sophistication of deepfake audio generation, detecting spoofed speech has become a critical challenge, particularly in multiparty conversations where multiple speakers interact dynamically. While research in the domain of audio spoof detection is primarily based on efficiently identifying individual audio clips as spoofed or not, this project will delve into scenarios where multiple speakers are involved. This system primarily aims to deal with observations reported by [5], that traditional audio spoof detection algorithms struggle to accurately predict audio streams, as spoofed or bonafide, in a multi-party environment. The system tries to determine the reason behind the drop in performance of spoof detection algorithms and possible solutions and/or strategies to help improve the detection accuracy. The proposed system will analyze multiparty conversations to identify if manipulated speech forms are present in a conversation, contributing to advances in the domain of audio deepfake detection.

# Chapter 2

# Background

This section explains the different basic concepts associated, based on which this system is built, and concepts that align closely with the system.

## 2.1 Methods of generation of spoofed audio

Spoofed audio refers to any kind of manipulated or fabricated speech designed to deceive a target victim or an automated system. Machine Learning algorithms could also be used to generate spoofed audio clips.

### 2.1.1 Replay Attacks

Replay attacks involve recording a person's voice, which is later used to impersonate the said person to deceive the target victim. Detection of these systems prove quite difficult as these audio clips do not undergo synthesis or modification, and thereby the underlying characteristics of the audio clips are maintained. Autonomous audio spoof detection systems rely on these modified characteristics to differentiate spoofed audio samples from bonafide samples. These characteristics will be discussed in later sections.

Replay attacks are of two different types: single hop and multi hop attacks. The single hop replay attack uses a direct playback of a recorded audio clip while a multi-hop replays audio through multiple devices. The characteristics of the audio clip degrade over several hops and, therefore, become more difficult to detect. ASVspoof 2019[6] compromises of single-hop audio clips and VSDC[7] proposes multi-hop replay attacks.

### 2.1.2 Speech Synthesis

Speech thesis technique generates artificial speech from text inputs, with the goal of producing realistic sounding audio that mimics a target speaker's voice. Generating speech

from text is referred to as Text-to-speech synthesis while modifying the voice's characteristics to resemble an attacker is called voice conversion. To accurately replicate the target speaker, the underlying text-to-speech ML model has to be trained on the person's speech samples to capture vocal nuances. This enables the ML model to accurately adjust the pitch, tone, speaking rate, and accent of the synthesized audio clip.

This technique often uses models such as Tacotron, which uses a Griffin-Lim algorithm over an RNN to generate raw waveforms, and a Multi-SpectroGAN that uses a conditional discriminator to eliminate reconstruction loss between ground truth and the generated mel-spectrogram [8]. Audio clips that have undergone similar syntheses form part of the [6] data set. Most spoof detection algorithms work on similar datasets to hone their prediction capabilities, to better ascertain a given input audio as either spoof or bonafide.



**Figure 2.1:** Common Spoofing Techniques [9]

### 2.1.3 Adversarial Attacks

To carry out such attacks, attackers introduce subtle changes to an audio waveform that are imperceptible to deceive human listeners. These modifications can significantly alter the output of deep learning detection models [10]. Attackers could have varying levels of knowledge about the underlying spoof detection model. If the attack has complete knowledge of the underlying ML algorithm, then the attack is termed a white-box attack. If the attacker has little knowledge, then it is a gray-box attack, while if the attacker has zero knowledge, then the attack is called a black-box attack [11].

**Figure 2.2:** Adversarial Attack [11]

Apart from differentiating adversarial attacks on the basis of the knowledge gathered about the underlying ML algorithm, they could also be classified in terms of techniques used:

- Signal Processing: Attackers could manipulate the different characteristics of an audio waveform such as the frequency spectrum or the phase of the audio signal. Time stretching and resampling can also be used to carry out this attack.
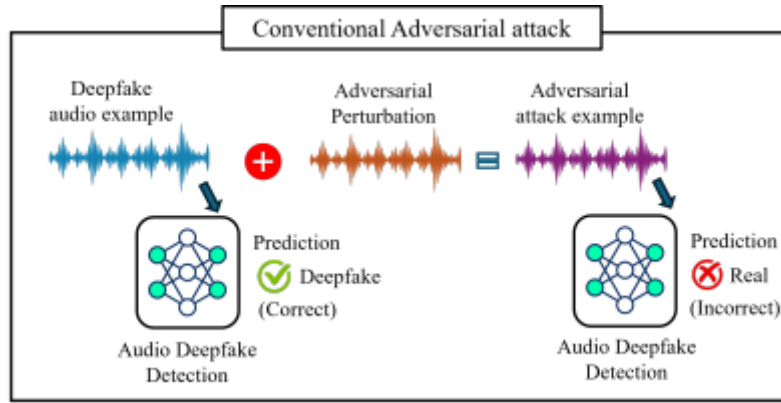
- Generative Adversarial Networks (GAN) attacks: in this attack, an audio generator uses a discriminator to align the spoofed sample more in line with the genuine sample [11]

## 2.2 Evaluation metrics

Evaluation metrics are standardized ways to assess how well an audio spoof detection algorithm performs. They help assess the reliability and accuracy of the system in question.

Equal Error Rate (EER) is the most commonly used evaluation metric in the domain of audio spoof detection. It represents a threshold where the value of the false rejection rate and the value of false acceptance rate coincide. In this context, the false acceptance rate represents the possibility of a spoof audio clip being classified as genuine. On the other hand, false rejection rate is the possibility that a genuine audio clip is falsely classified as spoofed. Spoof detection algorithms with a low EER perform better than their counterparts with a higher EER.

Accuracy signifies the percentage of audio samples correctly classified as being spoofed or genuine. Although accuracy is a common indicator of an ML algorithms' predictions, in the realm of audio spoof detection, EER tends to be a little more important than accuracy. This is the result of the datasets, used in the training of these algorithms, being heavily

unbalanced.

F1-score also has a role to play in an audio spoof detection algorithms' performance evaluation. It is the harmonic mean of precision and recall and is useful when dealing with datasets that are inherently biased, as it considers both the value of the false acceptance rate as well as the value of the false rejection rate. The difference between F1-score and EER is that the former has a fixed threshold, i.e. clear positive/negative labels assigned to predictions, whereas the latter deals with a variable threshold and is therefore better suited to deal with audio spoof detection.

Other commonly used evaluation metrics are precision and recall. Precision refers to how often the model's predictions of spoofed audio clips are actually spoofed. On the other hand, recall refers to how often actual spoofed clips are correctly classified as spoofed by audio spoof detection algorithms.

## 2.3  Audio Spoof Detection

As discussed above, Audio Spoof Detection is the process of classifying any given audio clip as spoofed or bonafide based on the features extracted from the said audio clip. In general, an audio spoof detection pipeline has four components, namely the Raw Audio Input, a Feature Extractor, A Classifier, and the prediction.

### 2.3.1  Raw Audio Input

The pipeline starts with a raw audio input, which can either be authentic or spoofed using techniques like text-to-speech, voice conversion or adversarial deepfake techniques as discussed previously. This input could either be in the form of a waveform or a digital audio file with extensions '.flac'.

### 2.3.2  Feature Extractor

Raw audio is fed into a feature extractor, where it is transformed into a meaningful representation that ML algorithms could utilize. These meaningful representations could be in the form of hand-crafted features or learned features. The former are designed by professionals based on established principles in audio signal processing and are usually based on properties such as frequency and/or energy. The latter is automatically learned by deep learning models, such as CNNs and Transformers, directly from the raw audio input. These models tend to be self-supervised. Some common hand-crafted features include:

1. Mel-Frequency Cepstral Coefficients (MFCC) - these help capture short-term frequency and timbral characteristics of speech, which can highlight the differences

between spoofed and genuine audio samples [12].

2. Constant Q Cepstral Coefficients (CQCC) - these provide a time-frequency break-down with pitch based spacing, making them useful in differentiating subtle cues that might reveal a spoofed artifact [12]

3. Spectrograms - Visual representations of the frequency spectrum over time, often used as visual input to CNNs to detect patterns or anomalies [12]

### 2.3.3  Classifier

These features that have been extracted above are then passed to a classifier which attempts to learn subtle differences between bonafide and spoofed audio samples. CNNs, ResNet, and Transformers are some of the most popular classifier architectures..

Convolutional Neural Networks (CNNs) are deep learning models that are designed to automatically learn spatial patterns in data. Although they were originally developed for image processing, they are well suited for audio tasks as well. These systems use convolutional filters that slide over the input such as a spectrogram, with each filter learning to detect specific patterns in the audio frequency-time space. While early layers capture simple patterns, deeper layers learn more complex structures making these systems robust to small variations and effective in extracting meaningful features from audio data [13].

Residual Network (ResNet) is a deep CNN architecture introduced by Microsoft that uses 'skip connections' to help train very deep networks. These connections allow the network to pass information across layers without modifications efficiently dealing with vanishing gradient problem. They have better feature extraction compared to CNNs, as they have higher depth. They also have better accuracy compared to standard CNNs [14].

Transformers are deep learning models originally designed for sequence data, but are now commonly used in audio tasks. They use a mechanism called self-attention to understand how different parts of the input relate to each other, no matter how far apart they are in the sequence. Unlike CNNs, transformers use layers of attention to evaluate the importance of each segment, making them effective at long-range patterns and contextual information in audio.
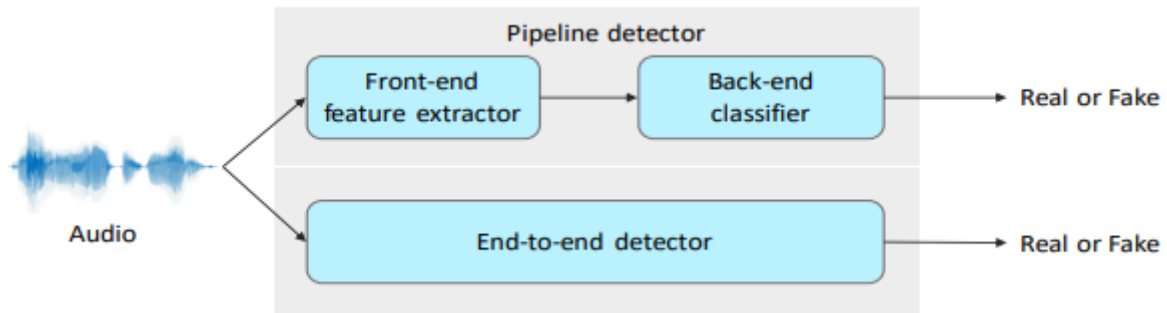
**Figure 2.3:** Spoof Detection Pipeline [15]

### 2.3.4  Prediction

The classifier finally outputs a prediction i.e. if the raw audio input is either spoofed or genuine. Typically, audio spoof detection is a binary classification problem. but could also provide confidence scores associated with the prediction or segment-level decisions for an in-depth analysis. This report presents audio spoof detection as a multi-class problem and will be explained in detail in the later sections.

## 2.4  Speech Separation

Speech separation is the process of isolating individual speakers' voices from a single audio recording where multiple people are talking at the same time. This is often called the 'cocktail party problem'. The system receives a mixed audio signal containing overlapping voices from two or more speakers. The system processes the input and aims to reconstruct the original speech signals for each individual speaker, producing separate audio channels for each.

Similar to Spoof Detection, the raw input undergoes pre-processing, where it is converted into a suitable representation. Most often, raw input is converted into waveforms or spectrograms. The former is a one-dimensional signal that shows how the air pressure changes over time i.e. how loud the sound is at each moment. It is a plot between the amplitude and time. The latter is a visual representation, as discussed above. The input is processed using the Short-Time Fourier Transform(STFT). The raw input is sliced into segments and each of these segments is subjected to STFT. It is a plot between time and frequency [16].

**Figure 2.4:** Spoof Detection Pipeline [16]

After the preprocessing is complete, the transformed input is fed into a feature extractor. In the case of transformed input being a spectrogram, the model learns to pick up patterns in time-frequency segments, which often carry speaker-speaker artifacts. In the case of transformed input being a waveform, the feature extractor works in a similar fashion by breaking down the waveform into a set of blocks.

Once the features are extracted, the deep learning models try to figure out which parts

of the transformed input belong to which speaker. Some techniques used include

- Mask Estimation - the model learns to create a 'mask' for each speaker. These masks indicate which segments of the input belong to which speaker. Applying these masks to the mixed audio allows the system to isolate individual speakers[17].

- Deep Clustering - the model projects time-frequency bins into an embedding space. In this space, bins that belong to the same speaker are grouped together. A clustering algorithm, then assigns each group to a specific speaker, thereby effectively differentiating the individual voices[18].

- Encoder-Separator-Decoder Architecture - these models employ an encoder which converts the raw waveform into a set of learnable segments. The separator module assigns speaker-specific weights to these components, after which the decoder reconstructs the waveform. This approach is commonly used in models such as Conv-TasNet, where a feature extractor step is unnecessary, as the encoder acts as a substitute within the larger separation architecture[18].

Unlike spoof detection models, speech separation models have a post-processing step. In this step, the separated features or masked signals are converted back into clean audio waveforms for each speaker. In the case of the use of spectrograms, the inverse STFT function is used over the segments to reconstruct the audio signal. On the other hand, for time-domain methods, the decoder reconstructs the waveform directly from the separated segments[18].

## 2.5   Transfer Learning

Transfer learning, a machine learning technique in which a model, initially trained to deal with a source task, is leveraged to deal with a target task. Traditionally, machine learning algorithms are built for each specific application, which would ideally have large amounts of associated data. They would also require substantial computational resources to deal with the size of the dataset. When a situation arises where an ML model has to be trained for a target task but lacks sufficient data or computational resources, Transfer Learning proves useful. A general Transfer Learning pipeline includes 3 steps:

1. Pre-training - In this step, a model is trained on a large, either a general purpose or a task-specific dataset. The model learns patterns and other parameters that could be used to solve a generic task[16].

2. Knowledge Transfer - The learned parameters i.e. the weights, from the pre-trained models are transferred to a new model designed for the target task[16].

3. Adaptation - The transferred model is either used as a feature extractor or fine-tuned using a smaller, task-specific dataset. While the former involves using learned

(a) Training a model for a different task without using transfer learning.

(b) Training a model for a different task using transfer learning.

**Figure 2.5:** Transfer Learning Usage[16]

representations as input for a new classifier or regressor, the latter allows some or all of the model's parameters to be fine-tuned to better fit the task at hand[16].

Transfer learning proves useful in scenarios where the availability of labeled data for a target task is limited. This is because this technique makes use of existing knowledge to reduce the time taken in training the model for the target task, and reduce the use of computational resources. This technique has resulted in increased accuracy and generalization[19] especially in computer vision, natural language processing, and medical imaging.

However, the effectiveness of this technique depends on the similarity between the original and target tasks. If the domains are too disjoint, the transferred knowledge may not be effective and can even reduce the performance of these models, commonly known as negative transfer[20]. Despite these challenges, transfer learning remains a useful tool for building useful ML pipelines with limited data resources.

# Chapter 3

# Related Works

This section compromises some of the most closely related scholarly articles and their findings have been summarized.

In one study [5], Arjun Pankajaskhan built a software tool designed to identify deepfake audio during live conversations across different communication platforms. Their approach tested ResNet and LCNN models on the ASVspoof 2019 dataset. Although the models performed well in static or pre-recorded environments, the tool struggled when dealing with the unpredictability of live, real-time calls like those on Microsoft Teams. In addition, they had voluntarily circumvented the 'cocktail party problem' by allotting time frames for each individual speaker in the call. The research highlights the gap between the performance of these models in static environments and in real-time dynamic environments.

Gustav A.P. Bonvang's work[16], focused on applying transfer learning with ResNet50 to detect audio deepfakes in real-world conditions. The model was trained on Mel spectrograms, derived from the In-The-Wild dataset. His results showed impressive gains, achieving almost 97% accuracy and 95% f1 score, outperforming non-transfer baselines by 22% accuracy and 44% f1 score respectively. The study illustrates how transfer learning can cut down training costs and data needs while maintaining strong detection results.

Govind Mittal et al.[21], developed an AI-powered verification system named PITCH, which uses a set of 20 audio-based challenges designed to test speaker authenticity. These include changes in tonal shifts, complexity of spoken phrases, and other variations. The system was tested on a massive dataset of deepfake audio and achieved an AUROC of 89% with machine-only detection, while human evaluators scored 72%. Combining human intuition with machine pre=screening increased the overall accuracy to 85%. The results support the idea that the combination of machine analysis with human judgment can strengthen voice authentication systems in practical scenarios.

A separate study [22], demonstrated how detectors can be manipulated using adversarial examples generated by GANs. These attacks were able to reduce the detection accuracy in both the digital and physical domains. The findings raise concerns about model vulnerability and call for strategies like adversarial training and input preprocessing(transformed representations are less sensitive to subtle, malicious changes) to defend against increasingly sophisticated spoofing methods.

He and Whitehill proposed a thorough survey of recent progress in end-to-end models for multi-speaker automatic speech recognition using single channel (monaural) audio[18]. The survey categorizes models into single-input-single-output and single-input-multiple-output paradigms, examining design trade-offs and recent architectural improvements. While the former deals better with an unknown number of speakers, it requires advanced decoding for complex overlaps, the latter requires fixed number of speakers, and they tend to work better under controlled settings. The study also emphasizes the differences between real-world datasets, such as LibriCSS, and simulated datasets, such as LibriMix. The former has better ecological validity but is more difficult to annotate accurately, while the latter enables better large-scale systematic experimentation but does not capture real-world scenarios accurately. The study highlights key challenges in this domain, such as overlapping speech, limited annotated data, and the difficulty of attributing words to individual speakers without special cues.

# Chapter 4

# Methodology

This Section provides a detailed explanation on how the system was formulated and adapted to the domain of identifying spoofed audio samples in a multi-party conversation.

## 4.1 Custom Dataset Creation

This project aims to develop a pipeline that is capable of identifying spoofed audio sources, in a multi-party environment which is often referred to as a cocktail party problem. As discussed in Section 2, the cocktail party problem refers to the challenge of distinguishing and focusing on individual audio sources, potentially spoofed in this context, where multiple speakers are present simultaneously. The primary challenge is the absence of relevant data. While existing resources such as the ASVspoof 2019 dataset consist of a mix of bonafide and spoofed audio samples, they are primarily designed for single-speaker verification tasks. They do not host audio samples from real or simulated multi-party conversations, where overlapping speech and/or interference is present. To address this limitation, this project proposes the creation of a novel dataset in an attempt to replicate the audio characteristics of a multi-party environment. This custom generated dataset includes mixtures of bonafide samples, spoofed samples, and an overlap of bonafide and spoofed audio samples.

To create a dataset of overlapping audio samples with desired characteristics for this project, a custom novel mixing process was implemented to generate 30,000 unique mixed files. The ASVspoof 2019 dataset was parsed, and the spoofed and bonafide audio samples were segregated into two arrays. The process of mixing begins by randomly selecting files which are then subsequently passed to the 'overlap_audio' function. To make the selection truly random, the cryptographically secure Python's 'secrets' module was used. This ensures the dataset isn't biased and is also resistant to patterns that might otherwise skew the learning efficiency.

The 'overlap_audio' function takes two input audio files- each labeled as either spoof or bonafide and generates a single output audio file. The overlapping technique induces realistic variations that could ideally be expected to take place in the presence of multiple speakers. The goal is to create a training dataset that better reflects the complexity of real-world multi-party conversational scenarios. The function performs the following steps:

1. During the pre-processing step, the audio files are loaded using the librosa package from python, preserving its original sampling rate. While loading, the audio clips are subjected to a check to ensure that the audio files use the same sampling rate, as mismatched rates can lead to distortion or misalignment during mixing.

2. A random delay is introduced to mimic the natural variability of the speaker's timings in audio calls. Here, one of the audio signals is randomly delayed by up to one second. This is done by padding the chosen audio with silence, introducing a forced delay.

3. Audio samples are then subject to length trimming, where the audio signals are trimmed to the same duration, to avoid mismatches when they are added together. The length of the shortened audio sample is chosen as the length of the resultant mixed audio sample.

4. Following this a gain adjustment is carried out where, each audio sample is scaled by a random factor between 0.7 and 1.3 to emulate variations in loudness between speakers. This difference in loudness could be the result of a difference in quality of microphone used, distance from their respective microphones, and the vocal intensity of the speakers.

5. The two audio samples are then mixed to produce a transformed waveform. This signal is also normalized to prevent clipping thereby maintaining consistent output levels.

6. Finally the output is labeled. A function creates a label for each resultant audio waveform. If the input audio samples are both bonafide, the output audio waveform is given a 'BB' label. Similarly, if both the audio samples are spoofed, the resultant audio waveform is given the 'SS' label. If a mixture of two is present, the output audio waveform is given a 'SB' label regardless of how the audios are mixed.

7. The resultant mixed audio is saved locally in the flac format, similar to the ASVspoof 2019 dataset, with a filename that includes the filenames of the individual audio samples.

The function outputs a mixed audio sample, along with metadata such as the filenames of the individual audio clips, their individual true labels, and the label of the combined audio waveform. This information is logged into a dictionary which is then used to create a single metadata file hosting all the information about all the audio samples generated.

## 4.2   Evaluation of SOTA spoof detection model on ASVspoof 2019

Before constructing a pipeline and running inference on our custom dataset, the state-of-the-art spoof detection model is tested on the original dataset. This is done to facilitate comparison between the results of this model on the individual audio samples and the audio samples that are a result of the speech separation process.

To perform spoof detection, we use a pre-trained WAVLM-based sequence classification model, fine-tuned on the ASVspoof 2019 dataset. The model, 'abhishtagatya/wavlm-base-960h-asv19-deepfake,' is loaded from the Hugging Face repository from the Hugging Face Transformers framework. The following components have been loaded:

- AutoConfig -this component loads the correct configuration class for a given pre-trained model. It contains information such as the number of transformer layers present, number of attention heads present, type and size of the convolutional feature encoder used, number of output classes, pooling method for aggregating sequence outputs, and dropout rates.

- Wav2Vec2FeatureExtractor - this component is used to convert raw audio waveforms into input features.

- WavLMForSequenceClassification - the component is the classification model itself. This component predicts whether the input audio sample is bonafide or if it has been spoofed.

```
1 config = AutoConfig.from_pretrained("abhishtagatya/wavlm-base-960h-asv19-
    deepfake")
2 feature_extractor = Wav2Vec2FeatureExtractor.from_pretrained("abhishtagatya/
    wavlm-base-960h-asv19-deepfake")
3 spoof_model = WavLMForSequenceClassification.from_pretrained("abhishtagatya/
    wavlm-base-960h-asv19-deepfake", config=config).to(device)
```

**Listing 4.1:** Loading the spoof detection model

A metadata file pertaining to the evaluation dataset of the ASVspoof 2019 dataset is parsed. For each audio file name that is parsed, an audio path, for the audio associated audio file, is constructed by appending '.flac' to the file name. This path is subsequently fed to the predict spoof function, which utilizes the in-built load_audio function to load the audio sample. The loaded raw waveform is passed to a SepFormer model, which separates the waveform into a two dimensional tensor which is split and stored into two individual arrays. Each of these arrays is passed to the 'Wav2vec2FeatureExtractor' essentially converting the audio signal into a transformed input which the model can process. The 'return_tensors="pt"' attribute is used to ensure that the output is a PyTorch tensor, while the 'Padding=True' attribute is used to handle variable length of the inputs by zero padding them. Finally, the 'torch.no_grad()' function is utilized to disable the gradient

tracking, ultimately leading to an increase in the efficiency of the inference process.

```python
# Part 2: Read Metadata from a .txt file

def load_metadata(metadata_file):
    metadata = []
    count = 0
    metadata_count = {}
    with open(metadata_file, "r") as f:
        for line in f:
            count += 1
            parts = line.strip().split(' ')
            filename, _, _, _, _, label = parts
            metadata.append((filename, label))
            if label not in metadata_count:
                metadata_count[label] = 1
            else:
                metadata_count[label] += 1

    print(count, metadata_count)
    return metadata

metadata = load_metadata(metadata_file)
print(metadata[:5])  # Display first 5 entries
```

**Listing 4.2:** Loading of custom metadata file

```python
def predict_spoof(audio_path):
    audio = load_audio(audio_path)
    inputs = feature_extractor(audio, sampling_rate=16000, return_tensors="pt", padding=True)
    with torch.no_grad():
        logits = spoof_model(**inputs.to(spoof_model.device)).logits
    probs = torch.softmax(logits, dim=-1).cpu().numpy()[0]
    predicted_label = "spoof" if probs[1] > probs[0] else "bonafide"
    return predicted_label, probs
```

**Listing 4.3:** Prediction function for spoof detection model

The model outputs raw logits for each class, which are subsequently converted into probabilities using the softmax function. Finally, the final labels are assigned based on which class has the higher probability. If the probabilities are equal, then the audio sample is classified as bonafide. Finally, the performance metrics accuracy, AUC score and EER of the inference is calculated for comparison purposes.

## 4.3 Evaluation of SOTA spoof detection model on the custom dataset

Following the construction of the overlapped audio samples, simulating the cocktail party problem, a testing pipeline is designed, combining a speech separation model and a spoof detection model to evaluate the feasibility of identifying spoofed audio sources in complex audio environments.

The system initally checks for GPU availability using Pytorch and assigns the computation to a CUDA-enabled GPU, for faster inference and better scalability when processing large batches of audio.

```
1 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
```
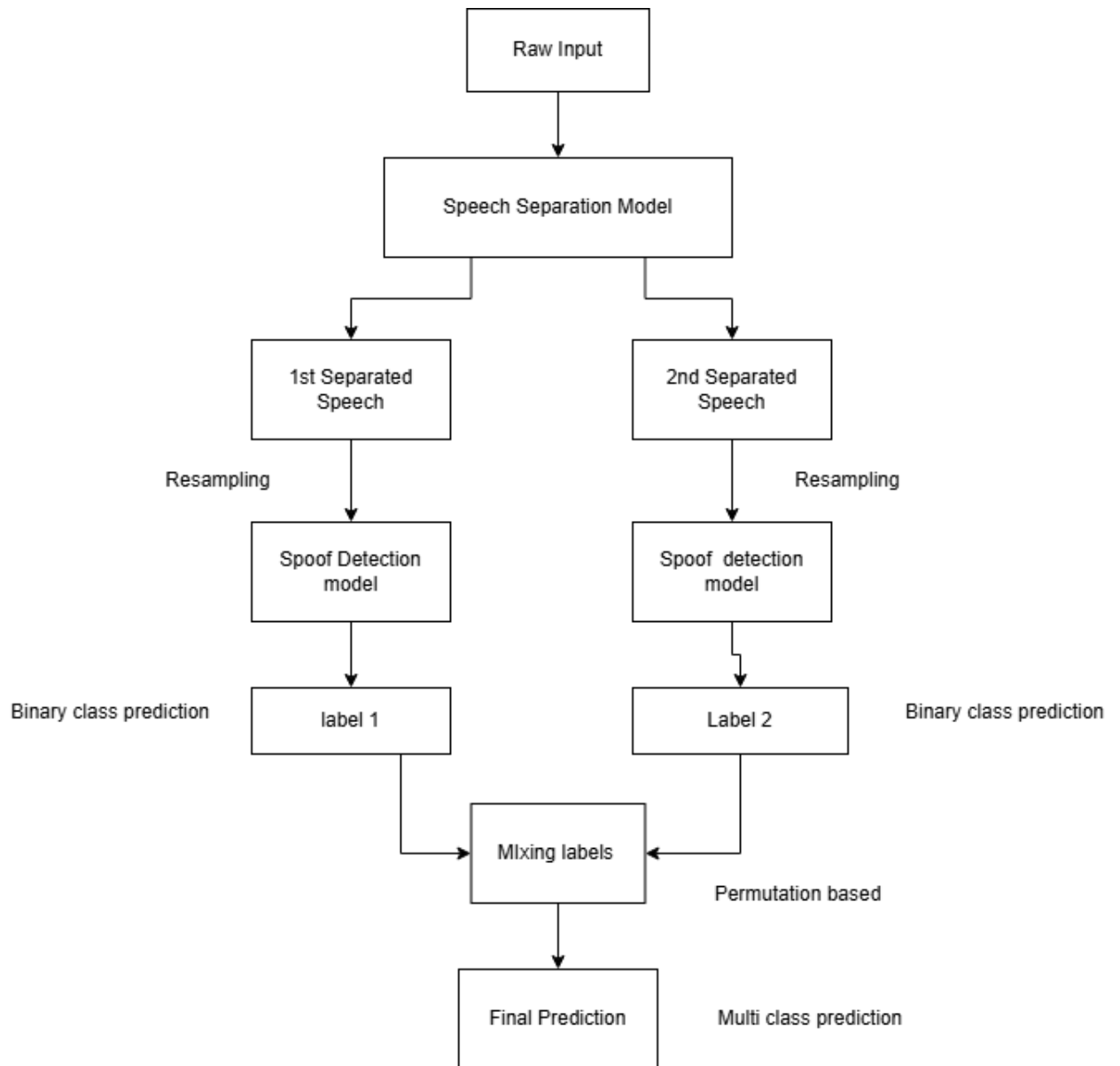**Listing 4.4:** Choosing CUDA-enabled GPU

To perform spoof detection, 'abhishtagatya/wavlm-base-960h-asv19-deepfake,' is loaded from the Hugging Face repository from the Hugging Face Transformers framework. The model has been loaded similar to the previous step. To handle overlapping speech, a state-of-the-art Speech separation model called the Sepformer has been incorporated from the Speech brain toolkit. This model is pre-trained on the WSJ0-2Mix dataset and has proven to provide high quality separation.

```
1 from speechbrain.pretrained import SepformerSeparation
2 model1 = SepformerSeparation.from_hparams("speechbrain/sepformer-wsj02mix")
```
**Listing 4.5:** Loading Speech Separation Model

The pipeline execution starts with the loading of the metadata files created in association with the mixed audio clips generation. The execution of the pipeline starts with a function iterating over the entire metadata file. The metadata file contains the mixed audio files' names and their associated true labels. The loop is set to 10,000 iterations as a result of a lack of memory and computational resources.

Each mixed audio sample is then passed through the previously loaded sepformer model which separates the mixtures into two speaker channels. The output is a tensor containing 2 clean speech waveforms. These waveforms are subsequently moved to the CPU, where they are individually resampled from 8Khz to 16Khz as the spoof detection model expects a audio sample at 16Khz. Each resampled audio waveform is then fed to a 'predict spoof' function, similar to the function described above, which outputs a label prediction i.e. spoof or bonafide and a probability score for each class. The individual predicted labels are converted into a string to resemble the true labels of the mixed audio clips. In other words, the output of the individual audio waveforms is combined to replicate the 'SS', 'BB', and 'SB' labels.

**Figure 4.1:** Architecture of the developed System

Since the order of the speakers in separation is not guaranteed, the pipeline introduces a permutation-invariant comparison step:

- it calculates the total prediction error using both the original and flipped order of speaker scores

- the permutataion with the lower error is chosen for the evaluation. This ensures a fair assessment regardless of which speaker was assigned to which output channel by the SepFormer.

This was performed to compare the predictions of the individual waveforms with their true values and not just the predictions of the combined audio sample. For efficient memory management and prevention of GPU overflow, temporary tensors and the CUDA memory cache are periodically deleted. At the end of the pipeline, the performance characteristics of the classification of mixed audio samples are calculated as well as the performance characteristics of spoof detection of individual audio waveforms post-speech separation. These evaluation metrics are calculated using the accuracy_score, precision_score, recall_score, and f1_score from the 'sklearn.metrics' package. The simplified architecture of the system developed has been illustrated in Figure 4.1.

## 4.4   Training of the Pipeline

In this section, the process of training the pipeline constructed in the previous section is discussed in detail. In the dataset creation section, it has been mentioned that only the evaluation dataset of the ASVspoof 2019 was used to create a custom evaluation dataset. For the purpose of training the pipeline, a custom development dataset and a custom training dataset were created from Development dataset and Training dataset from the ASVspoof 2019 dataset respectively. 30,000 mixed audio samples were created as part of the custom training dataset while 20,000 mixed audio samples were created as part of the custom training dataset. While the custom development training data was used as part of the validation part of an epoch, the custom training data was used to initially fine-tune the spoof detection model.

Following the dataset creation, the state-of-the-art spoof detection model and an equally competent speech separation model are used, similar to the process used in the above section. To mimic the training of the original spoof detection model, training parameters are used to train the spoof detection model as part of the pipeline. Here, the speech separation model has been frozen and only the spoof detection model has been targeted for fine-tuning, hence the use of the spoof model's learning parameters. The speech separation model has been frozen in order to facilitate the training of the spoof detection model when the other factors of the system are kept constant. This is to ensure that, the training solely focuses on adapting the spoof detection model to the characteristics of the mixed

dataset, preventing potential degradation of separation quality.

The AdamW optimizer, a variant of Adam that decouples weight decay from the gradient update has been used as it has consistently shown to improve the generalization capabilities of transformer-based models. The learning rate, Batch size and gradient accumulation are used in line with the training parameters of the original spoof detection model[23]. The parameters used are

- The learning rate is set to 1e-6 to prevent catastrophic forgetting and ensure gradual adaptation during fine-tuning.

- Batch size is set 1 due to memory constraints from large model size and input length

- Gradient accumulation is set to 2 so that gradients are accumulated over 2 steps before performing an optimizer step, to effectively simulate a larger batch size and stable updates

- Learning rate Scheduler is set to linear, therefore the learning rate will decrease linearly with increase in epoch number.

An epoch starts with the custom training data shuffled and sliced. Due to restrictions in computational availabilities, only 50% of the training data is used i.e. 15,000 custom training data is being used. The random seed of 42 is chosen for reproducibility reasons. The metadata is shuffled using this randomness, to prevent the loop from training on the same portion of the dataset, thereby eliminating redundant training. This increases the generalization capabilities of the pipeline.

```
1 random.seed(42)
2 random.shuffle(metadata_train)
3 metadata_train_subset = metadata_train[:15000]
```
**Listing 4.6:** Shuffling metadata array

When an audio sample passes through a sepformer, it undergoes separation and a resultant output with two isolated speakers is produced. The output is separated into two individual variables with each undergoing a squeeze function to ensure that the dimensions of the tensor and the input requirements of the spoof detection function match. Similar to the above section, the input waveform undergoes resampling and is converted into input features using the pre-trained Wav2vec2 feature extractor. The transformed input is fed to the spoof detection model with the true label (obtained from the metadata file) and the output loss is computed. Gradients are accumulated over multiple steps to simulate a large batch size, after which the following processes are carried out:

1. clipping of gradients to avoid exploding gradients.

2. modification of model weights by the optimizer

3. change in learning rate scheduler

4. resetting of gradients for the next batch

To prevent the system hosted on Google Colab from crashing, intermediate variables are deleted after the processing of each individual audio clip and CUDA memory is explicitly cleared every 100 files.

After each epoch, the model is evaluated on a held-out development set to measure the model's generalization capability during training. Due to restrictions in the availability of computational resources, a fixed-size subset of 10,000 audio samples is selected after the custom development metadata array is shuffled. This is to ensure an unbiased evaluation. Each audio file name in the metadata file is joined with an appropriate path that fetches the mixed audio sample from the custom development dataset. This audio sample is then fed into a SepFormer model which subsequently outputs two individual audio waveforms. Since the order of separated speakers is not guaranteed, all possible speaker-label alignments are evaluated. The permutation with lesser loss value is chosen as the final prediction. This ensures that label assignments aren't carried out random and/or mistakenly flipped. After determining the optimal permutation, the spoof detection model predicts labels for both audio samples, which are then compared to their true values. Classification losses are accumulated to compute the average loss per prediction.

Due to time constraints on how long a colab program could be run, the optimized parameters of the pipeline are saved, so they could be utilized in subsequent epochs. The code below is used to save the trained parameters. To enable fetching of the spoof model and feature_extractor components for training, the spoof model and feature_extractor components are saved. In addition, the current states of the optimizer and learning rate schedulers are saved to enable continuity in the training regime.

```python
from transformers import AutoConfig

save_path = "/content/drive/My Drive/My notes - cyber/Semester 4/
    finetuned_spoof_model6"

# Save model, feature extractor, and config
spoof_model.save_pretrained(save_path)
feature_extractor.save_pretrained(save_path)
spoof_model.config.save_pretrained(save_path)


# Save optimizer and scheduler
torch.save(optimizer.state_dict(), os.path.join(save_path, "optimizer.pt"))
```

```
13 torch.save(lr_scheduler.state_dict(), os.path.join(save_path, "scheduler.pt")
      )
14
15 print(f"Model, feature extractor, and config saved to {save_path}")
```

Listing 4.7: Saving learned parameters

Finally, to make the performance metrics of the epoch available, they are summed together as part of the metrics_history array and are then subsequently saved. This enables us to visualize the training metrics of the pipeline over several epochs.

```
1 metrics_history = {
2     "epoch": [],
3     "train_loss": [],
4     "val_loss": [],
5     "val_accuracy": []
6     # Add more if needed
7 }
8 metrics_history["epoch"].append(epoch + 1)
9 metrics_history["train_loss"].append(total_loss)
10 metrics_history["val_loss"].append(val_loss)
11 metrics_history["val_accuracy"].append(val_acc)
12 import os
13 import pandas as pd
14
15 metrics_file = "/content/drive/My Drive/My notes - cyber/Semester 4/
      training_metrics.csv"
16
17 # Check if file exists and load old metrics
18 if os.path.exists(metrics_file):
19     existing_df = pd.read_csv(metrics_file)
20     combined_df = pd.concat([existing_df, pd.DataFrame(metrics_history)],
      ignore_index=True)
21 else:
22     combined_df = pd.DataFrame(metrics_history)
23
24 # Save the updated file
25 combined_df.to_csv(metrics_file, index=False)
```

Listing 4.8: Saving performance metrics

For subsequent epochs, majority of the execution remains the same except for the loading of the spoof detection model and its parameters. Unlike the first epoch, where the model and its parameters are fetched from Hugging Face, the rest of the epochs utilize the parameters saved during the previous epochs. The parameters are loaded using the following code:

```
1 # Load optimizer and scheduler
2 optimizer.load_state_dict(torch.load(os.path.join(saved_model_path, "
      optimizer.pt")))
```

```
3 lr_scheduler.load_state_dict(torch.load(os.path.join(saved_model_path, "
      scheduler.pt")))
4
5 learning_rate = 3e-4
6 batch_size = 1
7 num_epochs = 4
8 gradient_accumulation_steps = 2
```

**Listing 4.9:** Loading saved paramters from previous epochs

The learning rate was increased in the final epoch alone, to speed up the learning process. This was a result of the lack of resources to carry out more epochs. Similar to the parameters, the model in itself was loaded from a locally saved instance.

```
1 # Check if CUDA is available for using GPU
2 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
3 print(f"Using device: {device}")
4
5 from speechbrain.pretrained import SepformerSeparation
6 model1 = SepformerSeparation.from_hparams(
7   "speechbrain/sepformer-wsj02mix"
8 )
9 saved_model_path = "/content/drive/My Drive/My notes - cyber/Semester 4/
      finetuned_spoof_model5"
10
11 config = AutoConfig.from_pretrained(saved_model_path)
12 feature_extractor = Wav2Vec2FeatureExtractor.from_pretrained(saved_model_path
      )
13 spoof_model = WavLMForSequenceClassification.from_pretrained(saved_model_path
      , config=config).to(device)
14 spoof_model.train()
15
16 spoof_model = spoof_model.to(device)
```

**Listing 4.10:** Loading saved model

## 4.5 Evaluation of the Final Trained Pipeline

Once training of the pipeline is complete, the trained pipeline is evaluated against the custom evaluation dataset. This is performed to calculate performance characteristics between the trained and pre-trained pipeline. To ensure an unbiased comparison, the same custom dataset has been used in both inference processes.

Although most of the pipeline characteristics were similar, the trained pipeline has small changes to how the models were loaded. While the speech separation model was loaded exactly the same way, it had been loaded during the inference of the pretrained

pipeline, the loading of the spoof detection model differed. Instead of the spoof detection model being loaded directly from Hugging face, the newly fine-tuned spoof detection model stored in the Google Drive was put to use. The code used to load this custom fine-tuned model is presented below:

```
# Check if CUDA is available for using GPU
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
print(f"Using device: {device}")

saved_model_path = "/content/drive/My Drive/My notes - cyber/Semester 4/
    finetuned_spoof_model6"

config = AutoConfig.from_pretrained(saved_model_path)
feature_extractor = Wav2Vec2FeatureExtractor.from_pretrained(saved_model_path
    )
spoof_model = WavLMForSequenceClassification.from_pretrained(saved_model_path
    , config=config).to(device)
spoof_model.eval()

spoof_model = spoof_model.to(device)
```

**Listing 4.11:** Loading final trained model

Similar to the pre-trained pipeline's execution, the system checks for the availability of a GPU. If a CUDA enabled GPU exists, all of the heavy computation is assigned to this device. The rest of the pipeline follows the traditional practice of loading the associated metadata file and iterating through it. Each file name encountered is appended to a local path to generate the associated local path of the mixed audio sample.

```
audio_path = os.path.join(audio_folder, filename)
```

**Listing 4.12:** Generating path for mixed audio files

This path is then used to fetch the mixed audio sample, which is fed to the SepFormer, effectively resulting in a two-dimensional tensor. Each tensor is unsqueezed, to resemble the input requirements of the spoof detection model, which transforms the same into processable units. These processable units are used by the spoof detection model to provide predictions of the individual audio sample. The predictions in the form of raw output logits are converted into human-readable probabilities, which are used to determine the predicted labels of the individual as well as the mixed audio sample. These labels are used to calculate the performance characteristics, as mentioned in the above subsection, for comparison purposes.

## 4.6 Ethical Considerations

Despite the role this spoof detection pipeline plays in multiparty spoofed scenarios, developing and testing such a pipeline also gives room for certain ethical issues that have to be addressed.

It is a prerequisite that the speech data, even if publicly available, must be governed by stringent privacy standards. This is to ensure that the unauthorized use of voice recordings, even in a synthetic or an anonymized form is reduced. This prevents possible infringement of privacy rights. All data sets used in this project originate from thoroughly vetted sources, who have also provided consent for use in research environments. This research does not collect any kind of personally identifiable information.

Automated systems have the tendency of inheriting biases from the data, they have been trained on. In the case of spoof detection models, the system could learn biases from the dataset, such as familiarity with certain accents, languages or speaker demographics. This could result in catastrophic outcomes with the system misclassifying certain groups more often. While this work, does not explicitly deal with these biases, maximum effort has been undertaken to create an even and unbiased custom dataset from the original ASVspoof 2019 dataset to avoid this misclassification problem.

As synthetic and deepfake speech technologies evolve, the boundaries between legitimate and malicious use blur. While building systems for the purpose of spoof detection isn't ethically wrong, the irresponsible deployment of them is. Often times, adversaries learn from existing spoof detection mechanisms, to further improve the spoofing techniques. To prevent such a scenario from occurring, the built pipeline has not been shared irresponsibly.

As a cyber security researcher, it is of utmost importance that the research is conducted ethically and it does not further endanger the existing digital space. The actions, mentioned above, were undertaken consciously to create a safer research environment and to ensure that the research does not, in any way, impede GDPR regulations or unwritten ethical rules.

# Chapter 5

# Results

The results obtained from the experiments in section chapter 4 is discussed in detail in this section.

## 5.1 Testing SOTA on the original dataset

Prior to the implementation of the spoof detection algorithm, the state-of-the-art spoof detection model was tested on the original ASVspoof 2019 for the sake of comparison between the model's performance across this dataset and the custom-generated evaluation dataset. The model had an impeccable accuracy at 0.9831, a precision of 0.9829, a recall of 0.9983, and an f1-score of 0.9905. The high accuracy indicates that the model correctly classified almost 99% of the data, showing a strong performance. The precision means that when the model predicts a sample as spoofed, then the sample actually belongs to the spoofed class. Higher recall means that the model minimizes false negatives. The confusion matrix of the spoof detection algorithm is visualized below:
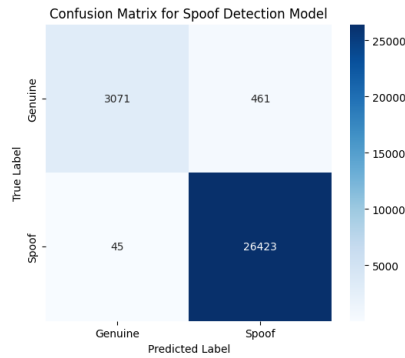


**Figure 5.1:** Confusion matrix of SOTA spoof detection model on ASVspoof 2019

The model has an EER of 0.0157 which is exceptionally good. The EER represents

the equilibrium point between the acceptance of false positives and false negatives. This low eer indicates that the model achieves a very low balance error between misclassifying spoofed audio samples as genuine or genuine samples as spoofed. This EER confirms the model's reliable performance in spoof detection.

## 5.2 Evaluation of Pre-trained Pipeline

To assess the performance of the pre-trained pipeline that incorporates the state-of-the-art spoof detection and speech separation models, we evaluated it on a custom-generated evaluation dataset. As discussed in Section chapter 4, this custom dataset compromises of overlapping speeches from, either two spoofed audio sources, two bonafide sources or a mix of one spoofed and one bonafide source. These overlapped audio samples were appropriately named 'SS,''SB,' and 'BB.' The classification system attempted to label each input mixed into one of these classes.

| Class | Precision | Recall | F1-score | Support |
|-------|-----------|--------|----------|---------|
| SS | 0.35 | 0.99 | 0.52 | 2535 |
| SB | 0.50 | 0.19 | 0.27 | 4979 |
| BB | 0.93 | 0.35 | 0.50 | 2486 |

**Table 5.1:** Classification metrics of the pre-trained pipeline

A class-wise analysis has been included below

- SS (spoof-spoof) - This class has a precision of only 0.35, thereby indicating that of all the samples predicted as SS, only 35% of them actually had SS as a true label. A recall of 0.99 suggests that nearly all audio samples of the SS class were correctly identified, hence the model is highly inclusive. This recall could be misleading as many of those predictions could be false positives. The f1 score suggests that while SS detection is inclusive, its precision is lacking.

- SB (spoof-bonafide) - This class has a precision of 0.50, so half of the samples predicted as SB were correct. The model has a recall of just 0.19, i.e. only 19% of the actual audio samples belonging to the SB class were correctly identified, indicating a major under-recognition of this class, i.e the model is selective. Finally its f1 score of 0.27 signifies difficulty in handling heterogenous cases.

- BB (bonafide-bonafide) - This class has a higher precision, compared to the other classes, at 0.93. This signifies that the model performs exceptional in avoiding false positives. But the model suffers with a 0.35 recall, which means that it misses most of the actual BB samples.

Looking at the bigger picture, the model has an overall accuracy of 42.98%. A random prediction for a model with 3 output classes is at 33.3%, so for the 10,000 samples, this

model performs just better than a model that randomly predicts. The overall weighted precision, weighted recall, and weighted F1-score of the model is a mere 0.59, 0.51 and 0.43 respectively. These values take into account the support variable, which indicates the number of instances for each class. The lower weighted f1 score suggests that the high misclassification of the more prevalent SB class significantly drags down the average.

To provide further insight, a confusion matrix was designed based on the classification performance across the three mixed-label categories. The confusion matrix further strengthens the inference that the model has a higher tendency to classify audio as SS. They also underscore the necessity of fine-tuning the model on the specific structure of the custom dataset.



**Figure 5.2:** Confusion matrix of pre-trained pipeline - multiclassification task

Apart from evaluating the model on the mixed label classification task, we assessed its performance on the individual audio streams as well. This evaluation reflects the models binary classification capabilities as the resultant audio stream could either be classified as spoofed or bonafide only. The model had an accuracy of 67.92%, a precision of 61.06%, a recall of 99.04%, and an f1 score of 75.62%. The high recall means that the model correctly identified nearly all spoofed audio samples, demonstrating extreme sensitivity towards cues associated with spoofed speeches. The precision is noticeably lower, indicating that a significant number of bonafide audio samples were misclassified as spoof. The f1 score balances the trade-off between the above two metrics, indicating that the model performs decently on the binary classification task. The Equal Error rate of the system was poor 23.90% at a threshold of 0.9999 suggesting that the model is heavily skewed toward avoiding spoof acceptance. Ideally, systems deployed for spoof detection have an EER of less than 10%. Therefore, this system is not suitable for deployment yet and requires further fine tuning.

**Figure 5.3:** Confusion of Pretained pipeline - binary classification

The model seems to correctly predict spoofed audio samples at the expense of misclassifying legitimate audio samples as spoof. Although most models in the high-stake audi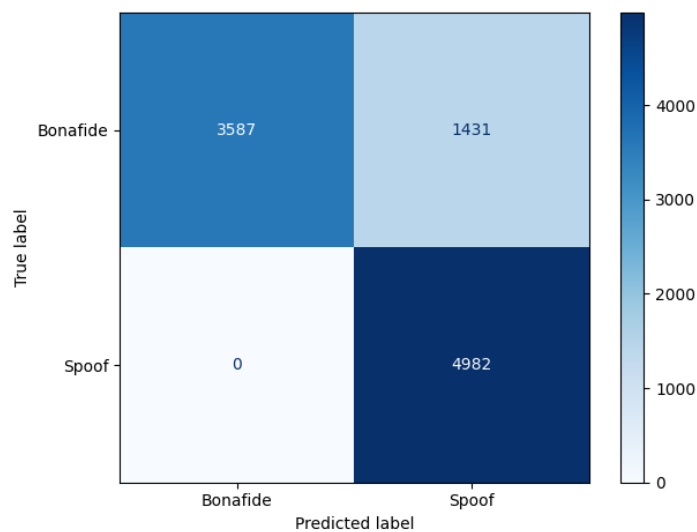o verification follow a conservative approach similar to this pipeline, the pipeline still seems to massively under-perform. The discrepancy in performance metrics between multi-class, and binary class classification indicates that while the spoof detection model performs average, an error is post-separation label assignment.

## 5.3 Train pipeline

In an attempt to increase the efficiency and the capability of the pipeline, the pipeline was trained on the custom generated mixed train and mixed development data. As discussed in Section chapter 4, the speech separation was frozen and only the parameters associated with the spoof detection model were fine tuned. For a total of 6 epochs, the models were trained over 15,000 mixed audio samples. The model was then evaluated over 10,000 mixed audio samples from the custom-generated development dataset.

The average training loss has evidently been decreasing across the epochs, indicating that the model has inferred patterns and is effectively learning and converging. The validation loss improved quickly during the first two epochs and plateaued after epoch the third epoch. This suggests that the model has learned generalizable patterns early on. In an attempt to push the model towards betterment, the learning rate was modified for the last epoch from 1e-6 to 3e-4. The validation loss increased infinitesimally but so did the validation accuracy. This analysis supports that the spoof detection model trained on separated audio streams was able to learn meaningful patterns that ultimately helped the

| Epoch | Cumulative Train Loss | Average Train Loss | Validation loss | Validation accuracy |
|-------|-----------------------|--------------------|-----------------|---------------------|
| 1     | 7731.96               | 0.5155             | 0.4292          | 0.8028              |
| 2     | 5657.80               | 0.3772             | 0.4028          | 0.8195              |
| 3     | 5120.00               | 0.3413             | 0.3836          | 0.8280              |
| 4     | 4750.49               | 0.3167             | 0.3826          | 0.8292              |
| 5     | 4626.12               | 0.3084             | 0.3826          | 0.8292              |
| 6     | 4915.02               | 0.3277             | 0.3831          | 0.8300              |

**Table 5.2:** Training metrics of the pipeline



**Figure 5.4:** Training loss vs validation loss

model generalize better on the validation data.

## 5.4   Trained pipeline

Following the horrendous results of the pre-trained pipeline, the spoof detection pipeline was trained. To be more precise, the spoof detection model was fine tuned while the speech separation model was frozen. This trained pipeline was subsequently evaluated on the custom dataset compromising of 10,000 mixed audio samples. The evaluation was performed similar to the how the pre-trained pipeline was evaluated.

Class wise analysis says that:

- SS (spoof-spoof) - the recall of the trained pipeline dropped from 99% to a mere 45% while the precision and f1 have increased slightly with the metrics being 55% and 49% respectively. This drop in recall suggests that the model is not skewed towards

| Class | Precision | Recall | F1-score | Support |
|-------|-----------|--------|----------|---------|
| SS | 0.55 | 0.45 | 0.49 | 2535 |
| SB | 0.63 | 0.38 | 0.47 | 4979 |
| BB | 0.47 | 0.93 | 0.62 | 2486 |

**Table 5.3:** Classification Report of Trained pipeline

flagging everything as spoofed and a balance in predictions across classes have been attained.

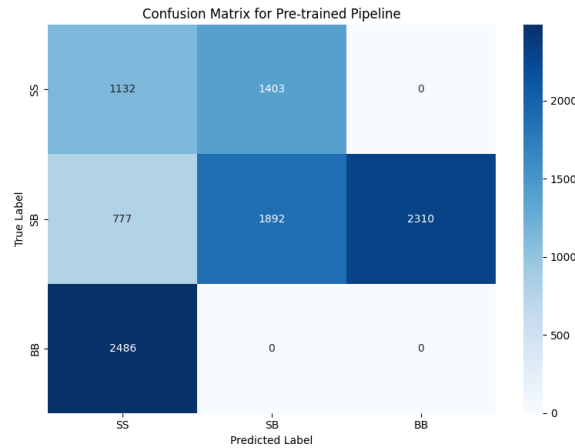- SB (spoof-bonafide) - The model attained an improvement in precision. The precision rose from 50% to 63% between pre-trained and trained pipelines. The recall also seems have to risen from 19% to 38%. The model's f1-score has also risen from 27% to a 47% indicating the model's fine tuning has proven useful. The model has learned better at predicting the instances from the SB class, which is more representative of the real-world cocktail party scenarios.

- BB (bonafide-bonafide) - The model had an increase across the recall and f1-score characteristics with the recall increasing from 35% to an impeccable 93% while the f1-score increasing from 50% to 62%. The fine-tuning of the model has improved it's ability to better recognize audio samples belonging to the BB class, and has also learned to avoid false positives on the spoof samples.

To provide further insights, a confusion matrix was formulated for the trained pipeline, similar to the pre-trained pipeline. The confusion matrix reasserts the previous inferences showing that the model shows a strong ability to correctly identify bonafide samples and that the SS class' precision and recall also seem to have improved.



**Figure 5.5:** Confusion matrix of trained pipeline - multiclassification

Apart from evaluating the model's performance on the mixed label classification task, the model's performance on individual audio streams retrieved after speech separation was also assessed. This analysis provides a clearer view of the underlying binary classification task in a different setting. The model had an accuracy of 74.41%, precision of 84.72%, recall of 59.87% and an F1-score of 70.16%. The model's high precision indicates that the model rarely raises a false alarm i.e. bonafide speech is likely to be misclassified as spoofed speech. The recall value has dropped miles in comparison to the pre-trained, indicating that the model's ability to identify all spoofed samples as spoof has reduced. An EER of 21.98% indicates that the model seems to have marginally better performance compared to the pre-trained pipeline. It indicates that at a threshold of 0.1935, the model misclassifies a spoofed sample as bonafide or vice versa, 21.98% of the time. The extremely low decision threshold implies that the model classifies inputs as spoof only if it is highly confident.
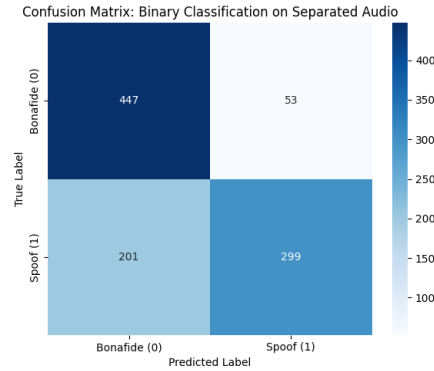


**Figure 5.6:** Confusion matrix of trained pipeline - binary classification

The discrepancy in performance metrics between multi class and binary class classification indicates that while the spoof detection model performs average, an error is post-separation label assignment is present.

# Chapter 6

# Discussion

This section discusses the impact of the research, the limitations faced during the study and the possible future works that can be done to capitalize on the results produced in this study.

## 6.1 Evaluation of SOTA spoof detection model on ASVspoof 2019 vs Custom dataset

In this study, it can be observed that the spoof's detection model's performance underwent a noticeable drop when transitioning from evaluating individual audio samples to mixed samples. When tested individually on the ASVspoof 2019 dataset, the model had an impeccable performance with an accuracy of 98% and an EER of 1.15%. These results clearly indicate the the model is only suited for clean, controlled and well-structured spoof detection tasks where the samples do not have any kind of interference from other speakers and acoustic condition. In comparison the spoof detection model had an accuracy of 42.98% over mixed audio samples and an accuracy of 67.92% on individual audio streams that were a result of mixed audio stream undergoing speech separation. In addition, the spoof detection model had an EER of 23.9%. These results signify that the model underwent a clear degradation in performance and the associated decision boundary sharpness compared to its counterpart. This degradation dropped further when the predictions were mixed to form the final mixed label prediction where the accuracy of the model dropped to a stooping 42.98%.

Most of these performance drops could be attributed to the speech separation process. When an audio sample undergoes speech separation, the temporal-spectral features of the audio sample are modified. The Zero-Crossing Rate(ZCR), which measures how often a signal crosses zero amplitude, is altered. Voiced speech tend to have a low ZCR, while Unvoiced Speech have high ZCR[24]. When an audio sample undergoes speech separation,

high ZCR regions are suppressed, which causes erasure of unvoiced spoofing artifacts. Certain spoof detection algorithms use Average ZCR as a cue for detecting spoofed audio samples[25]. So this erasure of key spoofing artifacts could have been the reason for the drop in performance of the spoof detection problem. In addition, it has been found that most queues for spoof detection lie in the range of 0 and 4kHz in normal voiced speech and in the range of 4-8khz in quiet parts or unvoiced sounds[26]. [27] quotes that their proposed speech separation method retains dominant peaks but the modulation effects in the frequency range of 50-200Hz are excluded which leads to a residual 'buzz' artifact. Since spoof detection also depends on these artifacts, the removal of the same affects its prediction abilities. Apart from these interferences, the audio sample also undergoes a resampling process as the audio sample that is output from the SepFormer is of a different sampling rate and therefore cannot be fed into the spoof detection model. This resampling process could have introduced artifacts that ultimately hinder the spoof detection model's capabilities.

As the ASVspoof 2019 dataset is designed specifically with spoofed samples containing spoofing artifacts that are a result of known spoofing techniques, a model trained on this dataset struggles to generalize on a mixed dataset with overlapped speech samples because of the elimination or inclusion of spoofing artifacts caused either by the mixing of speech samples and/or separation of mixed audio samples. The difference in results between the labeling of the overall mixed audio sample and the labeling of individual audio streams post-separation is a result of one of the speech samples being identified incorrectly, out of the two audio clips produced by the SepFormer. This inherently lowers the accuracy of the final pipeline.

## 6.2 Evaluation of Pre-trained vs Trained pipeline

To enhance the performance characteristics of the pre-trained model, the pipeline was fine-tuned to the custom-generated mixed dataset. The trained pipeline did witness an increase in performance as compared to its counterpart. This signifies that the training did help the model learn measurable and meaningful representations that ultimately boosted the pipeline's performance. While the overall gains appear moderate, they highlight important aspects of domain adaptation and robustness of speech-based ML systems.

While the f1 score for the SS class slightly declined for the trained pipeline in comparison to the pretrained pipeline, the f1 scores associated with the SB and BB class showed an improvement. The SB class showed the largest increase, indicating that the model learned key artifacts that help deal with inter-class variations introduced by the overlap algorithm. The binary classification also showed an improvement with the accuracy increasing by almost 7% in comparison to the pre-trained pipeline. There was an increase in precision but a decrease in recall, which ultimately led to a reduced EER.

This increase in performance could be attributed to the consistency of spoofing artifacts introduced in the speech separation process. As the speech separation model, SepFormer, was frozen, the separation artifacts introduced across the dataset would be consistent. The pipeline was able to grasp on to these consistent artifacts as noise in the classification task. These artifacts are considered noise, as they are predictable and are present across all the audio samples regardless of their true labels being spoof or bonafide. The possibility that the spoof model also adapted to the suppressed low-frequency artifacts is possible. The trained model could also recalibrate its ability to detect noise residuals that are a common modification to audio samples introduced by the speech separation process.

The reduction in EER by 2%, an increase in the binary classification's accuracy by 7% and an increase in multi classification's accuracy by 11% shows promising potential towards the development of a better system, which is better equipped in handling the cocktail party problem with spoofed audio samples.

## 6.3  Limitations

In the process of conducting this research, the system faced certain logistical and technical limitations. The primary challenge being the lack of time in furthering the research. Training over 15,000 data items and the model's validation over 10,000 data items over 6 epochs consumed ĩ00 hours. Besides the running time of the pipeline, analysis of the artifacts induced by speech separation that affect spoof detection is also heavily time consuming. Although the research indicates that these artifacts induce a noticeable performance reduction of the integrated spoof detection model, due to a lack of time, the research couldn't delve into which of these artifacts play a bigger role in performance reduction and by how much. Despite the use of A100 GPUs, an extremely powerful hardware, the training of the model took over 100 hours. In addition to these, inference on the evaluation dataset was also time consuming.

Due to hardware limitations and processing constraints, the batch size was reduced to just 1 which increases the risk of overfitting and may lead to poor generalization to unseen data. In addition, despite the generation of 30,000 training data and 20,000 validation data, only 15,000 training data and 10,000 validation were used during the training process. While 30,000 in itself isn't enough for the training process, further slicing of the dataset hindered the training process. Furthermore, the custom dataset generation was conducted solely based on the LA subset of the ASVspoof 2019 dataset. This could result in the pipeline overfitting to a specific type of spoofing attack, thus reducing it's generalization to other spoofing attacks. Due to hardware limitations, training of the entire pipeline was not conducted. In this research, the speech separation pipeline had to be frozen to save compute resources. While necessary, this technique has restricted the

end-to-end learning of the pipeline where the separation model could have worked hand-in-hand with the spoof detection model to improve the overall performance of the pipeline.

Another challenge faced by the pipeline is the adaptation of EER over a multi classification task. In the domain of spoof detection, EER plays a vital role in the evaluating the performance of the system. While most systems use EER over a binary classification task, this research attempted at adapting the EER over a multi classification task and failed. Although EER was still used in the system to evaluate its performance, EER was calculated based on the model's ability to ascertain a label to the audio samples that were derived from the speech separation task and not on the multi classification task of predicting if the audio sample belonged to the 'SS,''SB,' and 'BB.'

## 6.4   Future Works

This research while demonstrating the promise and limitation of the current setup, in an overlapped speech scenario, it also opens up opportunities for advancements in the fine-tuning of the system as well as certain other explorable avenues.

### 6.4.1   Fine-tuning the Speech Separation model

As discussed in the the previous section, the research presents with an opportunity to hone the speech separation model, in addition to the spoof detection model. In this research, the speech separation system was frozen to reduce complexity of the implemented system, requirement of computational resources and to isolate the fine-tuning of the spoof detection model. However, future studies could investigate end-to-end learning, where both modules are trained together simultaneously. This may allow the separation model to become spoof-aware, learning to retain discriminative artifacts necessary for spoof detection.

### 6.4.2   Increase in speakers and Integration of Speaker diarization

The current pipeline assumes only two speakers with binary spoof labels to be present in the custom generated dataset. Future work could be done on increasing the number of speakers involved in the cocktail party problem and investigating how an increase in speaker audio samples could affect the system's prediction capabilities. The system also works only on identifying the number of spoofed samples i.e. one or two spoofed samples. Future work could integrate a speech diarization model that would help the system attribute the labels to its specific speakers, which could prove useful in security contexts.

### 6.4.3  Dataset Generation

The system to its best capability avoided including inherent bias during dataset generation. Although cryptographic randomness was introduced to alleviate inclinations of dataset towards a specific class, the final dataset generated was not explored for the existence of possible prejudices. A promising direction is to investigate a better technology or process with recent and more diverse spoofing methods, including real-time recordings and AI-generated deepfakes to test the robustness of the system in more adversarial conditions. The custom generated mixed audio samples also lacks replication of real-world acoustic variations. Real-world variations such as reverberation, environmental noise, inconsistencies with the microphone in use or adversarial attacks could provide additional insights that the system could adapt to, making it relevant to the requirements of the deployable system.

### 6.4.4  Latency of the pipeline

The goal of the project was limited to just exploring the relevance of state-of-the-art spoof detection algorithms, in a overlapped audio conversation scenario which are representative of group calls in platforms such as Teams and Discord, and podcasts that could be subject to spoofing attacks. The system does not delve into increasing the efficiency of the training and inference processes. The system clearly required excessive time periods during the training and inference process, which makes it not so suitable for real-time spoof detection deployments. Future work could explore methods and techniques that help optimize the inference algorithm and provide low-latency models that are better suited for real-world situations.

# Chapter 7

# Conclusion

This chapter summarizes the experiments carried out in this research, the scope of the research, and the results produced by the research.

The research began with the goal of analyzing the performance of spoof detection model, which was trained on the simple ASVspoof 2019 dataset with individual audio samples, across a custom generated dataset which was more representative of a cocktail party problem with spoofed audio samples. The spoof detection model was initially tested on the ASVspoof 2019 dataset. The model performed exceptionally well with an accuracy of 98% and an EER of 1.15%. The model was later integrated into a pipeline with a speech separation model, SepFormer, which helped split the mixed audio samples. The resultant audio clips were tested for the spoofing artifacts by the spoof detection model. The untrained pipeline had an accuracy of 43% and an EER of 23%. In an attempt, the pipeline was fine-tuned on the custom generated dataset, with the speech separation model kept frozen. The fine-tuned pipeline was evaluated on the custom-generated evaluation dataset and the model returned an accuracy of 53%.

The research clearly indicates a drop in the spoof detection model's performance as a result of the introduction and erasure of key artifacts by the speech separation model. The report also suggests that further work could be done to identify which of these artifacts play a more important role in the performance of the pipeline. The report also emphasizes the need for additional time and hardware resources for furthering the research in this domain. The report also necessitates the need for a speech diarization model that could attribute which of the speakers in the conversation are spoofed.

The report concludes that, although spoof detection models with excellent prediction capabilities exist, these models struggle in a scenario that includes multiple speakers. This task requires the need for a domain-specific pipeline, such as the one implemented in the scenario, to deal better with the cocktail party problem with possible spoofed speakers.

# Bibliography

[1] Zahra Khanjani, Gabrielle Watson, and Vandana P. Janeja. "Audio deepfakes: A survey". In: *Frontiers in Big Data* Volume 5 - 2022 (2023). ISSN: 2624-909X. DOI: 10 . 3389/fdata.2022.1001063. URL: https://www.frontiersin.org/journals/big-data/articles/10.3389/fdata.2022.1001063.

[2] Mehmet Karakose et al. "A New Approach for Deepfake Detection with the Choquet Fuzzy Integral". eng. In: *Applied sciences* 14.16 (2024), pp. 7216–. ISSN: 2076-3417.

[3] April Rubin. *Scams use AI to mimic senior officials' voices, FBI warms.* https://www.axios.com/2025/05/15/artificial-intelligence-voice-scams-government-officials?utm_source=chatgpt.com. 2025.

[4] Times of India. *Scammers use AI George Clooney to steal over Rs. 11 lakhs from Facebook user.* https://timesofindia.indiatimes.com/technology/artificial-intelligence/scammers-use-ai-george-clooney-to-steal-over-rs-11-lakhs-from-facebook-user/articleshow/121171381.cms?utm_source=chatgpt.com. 2025.

[5] Jonat John Mathew et al. "Towards the Development of a Real-Time Deepfake Audio Detection System in Communication Platforms". eng. In: (2024).

[6] Xin Wang et al. "ASVspoof 2019: A large-scale public database of synthesized, converted and replayed speech". eng. In: *Computer speech language* 64 (2020), pp. 101114–. ISSN: 0885-2308.

[7] "Voice spoofing detection corpus for single and multi-order audio replays". eng. In: *Computer speech language* 65 (2021), pp. 101132–. ISSN: 0885-2308.

[8] Rishabh Ranjan, Mayank Vatsa, and Richa Singh. "Uncovering the Deceptions: An Analysis on Audio Spoofing Detection and Future Prospects". eng. In: (2023).

[9] Bowen Zhang et al. "Audio Deepfake Detection: What Has Been Achieved and What Lies Ahead". eng. In: *Sensors (Basel, Switzerland)* 25.7 (2025), pp. 1989–. ISSN: 1424-8220.

[10] Piotr Kawa, Marcin Plata, and Piotr Syga. "Defense Against Adversarial Attacks on Audio DeepFake Detection". eng. In: (2022).

[11]   Muhammad Umar Farooq et al. "Transferable Adversarial Attacks on Audio Deep-
       fake Detection". eng. In: (2025).

[12]   Taiba Majid Wani et al. "Detecting Audio Deepfakes: Integrating CNN and BiLSTM
       with Multi-Feature Concatenation". eng. In: *Proceedings of the 2024 ACM Workshop on
       Information Hiding and Multimedia Security*. New York, NY, USA: ACM, 2024, pp. 271–
       276. ISBN: 9798400706370.

[13]   Shuhui Qu et al. "Understanding Audio Pattern Using Convolutional Neural Net-
       work From Raw Waveforms". eng. In: (2016).

[14]   Noussaiba Djeffal et al. "Transfer Learning-Based Deep Residual Learning for Speech
       Recognition in Clean and Noisy Environments". eng. In: (2025).

[15]   Jiangyan Yi et al. "Audio Deepfake Detection: A Survey". eng. In: (2023).

[16]   Gustav Arnt Palmelund Bonvang. *Listening Beyond Words: Transfer Learning for Audio
       Deepfake Detection*. eng. 2023.

[17]   Peter Ochieng. "Speech Separation Based on Pre-trained Model and Deep Modu-
       larization". In: *ICASSP 2025 - 2025 IEEE International Conference on Acoustics, Speech
       and Signal Processing (ICASSP)*. 2025, pp. 1–5. DOI: 10.1109/ICASSP49660.2025.
       10888747.

[18]   Xinlu He and Jacob Whitehill. *Survey of End-to-End Multi-Speaker Automatic Speech
       Recognition for Monaural Audio*. 2025. arXiv: 2505.10975 [cs.CL]. URL: https://
       arxiv.org/abs/2505.10975.

[19]   Javier Canales Luna. *What is Transfer Learning in AI? An Introductory Guide with Ex-
       amples*. https://www.datacamp.com/blog/what-is-transfer-learning-in-ai-an-
       introductory-guide. 2024.

[20]   Karl Weiss, Taghi M Khoshgoftaar, and DingDing Wang. "A survey of transfer learn-
       ing". In: *Journal of Big data* 3 (2016), pp. 1–40.

[21]   Govind Mittal et al. "PITCH: AI-assisted Tagging of Deepfake Audio Calls using
       Challenge-Response". eng. In: (2024).

[22]   Mouna Rabhi, Spiridon Bakiras, and Roberto Di Pietro. "Audio-deepfake detection:
       Adversarial attacks and countermeasures". eng. In: *Expert systems with applications*
       250 (2024), pp. 123941–. ISSN: 0957-4174.

[23]   abhishtagatya. *abhishtagatya/wavlm-base-960h-asv19-deepfake*. https://huggingface.
       co/abhishtagatya/wavlm-base-960h-asv19-deepfake/blob/main/README.md?
       code=true. 2024.

[24]   Rafizah Mohd Hanifa et al. "Voiced and unvoiced separation in malay speech using
       zero crossing rate and energy". eng. In: *Indonesian Journal of Electrical Engineering and
       Computer Science* 16.2 (2019), pp. 775–. ISSN: 2502-4752.

[25]  Jinyu Zhan et al. "Detecting Spoofed Speeches via Segment-Based Word CQCC and Average ZCR for Embedded Systems". In: *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 41.11 (2022), pp. 3862–3873. DOI: 10.1109/TCAD.2022.3197531.

[26]  Menglu Li, Yasaman Ahmadiadli, and Xiao-Ping Zhang. "Audio Anti-Spoofing Detection: A Survey". eng. In: (2024).

[27]  Pejman Mowlaee, Mads Græsbøll Christensen, and Søren Holdt Jensen. *Sinusoidal masks for single channel speech separation*. eng. 2010.

# Appendix A

# Appendix

```python
1  # -*- coding: utf-8 -*-
2  """Overlap generation.ipynb
3
4  Automatically generated by Colab.
5
6  Original file is located at
7      https://colab.research.google.com/drive/1jOeyaZrMK0UAE2LkPO-fL9haJb05cTgQ
8  """
9
10 from google.colab import drive
11 import os
12 import librosa
13 import numpy as np
14 import soundfile as sf
15 import secrets  # Cryptographic randomness
16 import random  # General randomness
17 import torchaudio
18 import torchaudio.transforms as T
19
20 drive.mount('/content/drive')
21
22 # Step 2: Define dataset paths
23 base_path = "/content/drive/My Drive/My notes - cyber/Semester 4"
24 metadata_path = os.path.join(base_path, "ASVspoof2019_LA_cm_protocols")
25 output_path = os.path.join(base_path, "overlapped_audio_2")
26 metadata_output_path = os.path.join(base_path, "overlapped_audio_metadata_2.
       txt")
27
28 # Ensure output directory exists
29 os.makedirs(output_path, exist_ok=True)
30
31 # Step 3: Read metadata to classify audio files
32 spoof_files = []
33 bonafide_files = []
```

```
34 file_status = {}  # Dictionary to store file name -> bonafide/spoof
35
36 # Read the metadata file
37 metadata_file = os.path.join(metadata_path, "ASVspoof2019.LA.cm.eval.trl.txt"
      )
38
39 with open(f"{metadata_file}", "r") as f:
40     for line in f:
41         parts = line.strip().split()
42         if len(parts) < 5:
43             continue  # Skip malformed lines
44         _, filename, _, _, label = parts  # Extract relevant columns
45         file_path = os.path.join(base_path, "ASVspoof2019_LA_eval/flac",
      filename+'.flac')  # Adjust path
46         if os.path.exists(file_path):  # Ensure file exists before adding
47             if label == "bonafide":
48                 bonafide_files.append(file_path)
49             else:
50                 spoof_files.append(file_path)
51
52
53 print(f"Loaded {len(bonafide_files)} bonafide and {len(spoof_files)} spoof
      files.")
54
55 import torch
56
57 # Step 4: Function to create a realistic overlap and log metadata
58 def overlap_audio(file1, file2, label1, label2, output_dir):
59     audio1, sr1 = librosa.load(file1, sr=None)
60     audio2, sr2 = librosa.load(file2, sr=None)
61
62     if sr1 != sr2:
63         raise ValueError("Sampling rates do not match!")
64
65     # Generate a random time offset (up to 1 second)
66     max_delay = int(sr1 * np.random.uniform(0, 1))  # Up to 1 second shift
67     if secrets.randbelow(2) == 0:  # Randomly decide who starts first
68         audio1 = np.pad(audio1, (max_delay, 0))  # Delay first audio
69     else:
70         audio2 = np.pad(audio2, (max_delay, 0))  # Delay second audio
71
72     # Match length of both signals (cut longer one)
73     min_len = min(len(audio1), len(audio2))
74     audio1 = audio1[:min_len]
75     audio2 = audio2[:min_len]
76
77     # Apply random gain variation (loudness change)
78     gain1 = np.random.uniform(0.7, 1.3)  # Random loudness for audio1
79     gain2 = np.random.uniform(0.7, 1.3)  # Random loudness for audio2
80     audio1 *= gain1
81     audio2 *= gain2
```

```python
82
83      # # Apply reverberation (simulates room acoustics)
84      # reverb = T.Reverberate()
85      # audio1 = reverb(torch.tensor(audio1).unsqueeze(0)).squeeze(0).numpy()
86      # audio2 = reverb(torch.tensor(audio2).unsqueeze(0)).squeeze(0).numpy()
87
88      # Mix both audio signals
89      mixed_audio = audio1 + audio2
90
91      # Normalize to prevent clipping
92      mixed_audio = mixed_audio / np.max(np.abs(mixed_audio))
93
94      # Define mixed audio status
95      if label1 == "spoof" and label2 == "spoof":
96          mixed_label = "SS"  # Spoof-Spoof
97      elif label1 == "bonafide" and label2 == "bonafide":
98          mixed_label = "BB"  # Bonafide-Bonafide
99      else:
100         mixed_label = "SB"  # Spoof-Bonafide
101
102     # Save the mixed audio file
103     new_filename = f"{os.path.basename(file1).split('.')[0]}_{os.path.
    basename(file2).split('.')[0]}.flac"
104     output_filepath = os.path.join(output_dir, new_filename)
105     sf.write(output_filepath, mixed_audio, sr1)
106
107     return new_filename, file1, label1, file2, label2, mixed_label
108
109 # Step 5: Use cryptographic randomness to create 25,000 mixed audio files and
        log metadata
110 mix_log = []
111 num_pairs = 30000  # Target number of mixed audio files
112
113 for i in range(num_pairs):
114     if i % 5000==0:
115       print(i)
116     # Cryptographic randomness for true unpredictability
117     if secrets.randbelow(2) == 0:  # 50% chance of spoof-bonafide mix
118         file1 = secrets.choice(spoof_files)
119         file2 = secrets.choice(bonafide_files)
120         label1 = "spoof"
121         label2 = "bonafide"
122     else:
123         # Randomly choose bonafide-bonafide or spoof-spoof
124         if secrets.randbelow(2) == 0 and len(bonafide_files) > 1:
125             file1, file2 = random.sample(bonafide_files, 2)
126             label1 = label2 = "bonafide"
127         else:
128             file1, file2 = random.sample(spoof_files, 2)
129             label1 = label2 = "spoof"
130
```

```
131      # Overlap the chosen files
132      mixed_file, file1_name, label1, file2_name, label2, mixed_label =
         overlap_audio(file1, file2, label1, label2, output_path)
133      mix_log.append(f"{mixed_file} {os.path.basename(file1_name)} {label1} {os
         .path.basename(file2_name)} {label2} {mixed_label}")
134
135
136 # Step 6: Save the metadata for the new dataset
137 with open(metadata_output_path, "w") as f:
138     for entry in mix_log:
139         f.write(entry + "\n")
140
141 print(f"    Successfully created {len(mix_log)} mixed audio files!")
142 print(f"      Metadata saved at: {metadata_output_path}")
143
144 counter = {}
145 for line in mix_log:
146   parts = line.strip().split(' ')
147   if parts[-1] not in counter:
148     counter[parts[-1]]=0
149   counter[parts[-1]]+=1
150
151 print(counter)
```

**Listing A.1:** Complete code for Overlap Generation

```
1 # -*- coding: utf-8 -*-
2 """Pretrained sota spoof.ipynb
3
4 Automatically generated by Colab.
5
6 Original file is located at
7     https://colab.research.google.com/drive/1wtqimvr7pbYLldir1ucYycy6MR5-3iDH
8 """
9
10 !pip install torch torchaudio transformers scikit-learn soundfile asteroid
11
12 # Part 2: Load pretrained models and dependencies
13 import torch
14 from asteroid.models import ConvTasNet
15 from transformers import AutoProcessor, AutoModelForPreTraining,
       AutoModelForAudioClassification, AutoConfig
16 from transformers import Wav2Vec2Processor, Wav2Vec2ForSequenceClassification
       , Wav2Vec2FeatureExtractor, WavLMForSequenceClassification
17 import soundfile as sf
18 import os
19 import torchaudio
20 import torch.nn as nn
21 import torch.nn.functional as F
22 import matplotlib.pyplot as plt
23 from sklearn.metrics import roc_curve, auc
```

```python
24 import numpy as np
25 from collections import Counter
26
27 import torchaudio.transforms as T
28
29 # Part 1: Mount Google Drive
30 from google.colab import drive
31 drive.mount('/content/drive', force_remount=True)
32
33 # Set paths for the audio and metadata directories
34 audio_folder = "/content/drive/My Drive/My notes - cyber/Semester 4/
      ASVspoof2019_LA_eval/flac"
35 metadata_file = "/content/drive/My Drive/My notes - cyber/Semester 4/
      ASVspoof2019_LA_cm_protocols/ASVspoof2019.LA.cm.eval.trl.txt"
36
37 # Check if CUDA is available for using GPU
38 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
39 print(f"Using device: {device}")
40
41 # Load the pretrained separation model (ConvTasNet)
42 sepformer = ConvTasNet.from_pretrained('mpariente/ConvTasNet_WHAM!_sepclean')
43 sepformer = sepformer.to(device)
44
45 # config = AutoConfig.from_pretrained("abhishtagatya/wavlm-base-960h-itw-
      deepfake")
46 # feature_extractor = Wav2Vec2FeatureExtractor.from_pretrained("abhishtagatya
      /wavlm-base-960h-itw-deepfake")
47 # spoof_model = WavLMForSequenceClassification.from_pretrained("abhishtagatya
      /wavlm-base-960h-itw-deepfake", config=config).to(device)
48
49
50 config = AutoConfig.from_pretrained("abhishtagatya/wavlm-base-960h-asv19-
      deepfake")
51 feature_extractor = Wav2Vec2FeatureExtractor.from_pretrained("abhishtagatya/
      wavlm-base-960h-asv19-deepfake")
52 spoof_model = WavLMForSequenceClassification.from_pretrained("abhishtagatya/
      wavlm-base-960h-asv19-deepfake", config=config).to(device)
53
54
55 spoof_model = spoof_model.to(device)
56
57 from asteroid.models import BaseModel
58
59 model = BaseModel.from_pretrained("JorisCos/ConvTasNet_Libri2Mix_sepnoisy_16k
      ")
60 model = model.to(device)
61
62 !pip install speechbrain
63
64 from speechbrain.inference.separation import SepformerSeparation as separator
65
```

```
66 model1 = separator.from_hparams(source="speechbrain/sepformer-wsj02mix",
      savedir='pretrained_models/sepformer-wsj02mix',run_opts={"device": "cuda"
      })
67 model1 = model1.to(device)
68
69 import pandas as pd
70
71 # Load file, specify only needed columns
72 df = pd.read_csv(metadata_file, header=None, delim_whitespace=True, usecols
      =[1, 4], names=["file", "label"])
73
74 # Map labels to numeric (spoof = 1, bonafide = 0)
75 label_map = {"spoof": 1, "bonafide": 0}
76 df["label"] = df["label"].map(label_map)
77
78 df.shape
79
80 def load_audio(audio_path):
81     waveform, sample_rate = torchaudio.load(audio_path)
82     if sample_rate != 16000:
83         waveform = torchaudio.transforms.Resample(sample_rate, 16000)(
      waveform)
84     return waveform.squeeze().numpy()
85
86 def predict_spoof(audio_path):
87     audio = load_audio(audio_path)
88     inputs = feature_extractor(audio, sampling_rate=16000, return_tensors="pt
      ", padding=True)
89     with torch.no_grad():
90         logits = spoof_model(**inputs.to(spoof_model.device)).logits
91     probs = torch.softmax(logits, dim=-1).cpu().numpy()[0]
92     predicted_label = "spoof" if probs[1] > probs[0] else "bonafide"
93     return predicted_label, probs
94
95 from tqdm import tqdm
96 from sklearn.metrics import accuracy_score, confusion_matrix, roc_auc_score
97
98 y_true = []
99 y_pred = []
100 y_score = []
101 count = 0
102
103 for _, row in tqdm(df.iterrows(), total=len(df)):
104     count += 1
105     if count > 30000:
106         break
107
108
109     audio_path = os.path.join(audio_folder, row["file"])
110     audio_path = audio_path + '.flac'
111
```

```
112    if not os.path.exists(audio_path):
113        print(f"Missing: {audio_path}")
114        continue
115
116    pred, score = predict_spoof(audio_path)
117
118    y_true.append(row["label"])
119    y_pred.append(pred)
120    y_score.append(score)
121
122 print(y_true)
123 print(y_pred)
124 print(y_score)
125
126 # === Step 6: Evaluate ===
127
128 y_true_label = ["spoof" if label == 1 else 0 for label in y_true]
129 cm = confusion_matrix(y_true_label, y_pred)
130
131
132 acc = accuracy_score(y_true_label, y_pred)
133
134 y_score_modified = [score[1] for score in y_score]
135 roc = roc_auc_score(y_true, y_score_modified)
136
137 print(f"\nAccuracy: {acc:.4f}")
138 print("Confusion Matrix:\n", cm)
139 print(f"ROC AUC Score: {roc:.4f}")
140
141 from sklearn.metrics import roc_curve
142 import numpy as np
143
144 fpr, tpr, thresholds = roc_curve(y_true, y_score_modified, pos_label=1)
145
146 # Find the point where FPR    1 - TPR
147 eer = fpr[np.nanargmin(np.absolute((1 - tpr) - fpr))]
148
149 print(f"Equal Error Rate (EER): {eer:.4f}")
```

**Listing A.2:** Complete code for Testing SOTA spoof detection model

```
1  # -*- coding: utf-8 -*-
2  """pretrained pipeline.ipynb
3
4  Automatically generated by Colab.
5
6  Original file is located at
7      https://colab.research.google.com/drive/1HveSqt1YMWgGNFNq05qH9QpwqR8YXB9w
8  """
9
10 !pip install torch torchaudio transformers scikit-learn soundfile asteroid
```

```
11
12 !pip install speechbrain
13
14 # Part 2: Load pretrained models and dependencies
15 import torch
16 from asteroid.models import ConvTasNet
17 from transformers import AutoProcessor, AutoModelForPreTraining,
     AutoModelForAudioClassification, AutoConfig
18 from transformers import Wav2Vec2Processor, Wav2Vec2ForSequenceClassification
     , Wav2Vec2FeatureExtractor, WavLMForSequenceClassification
19 import soundfile as sf
20 import os
21 import torchaudio
22 import torch.nn as nn
23 import torch.nn.functional as F
24 import matplotlib.pyplot as plt
25 from sklearn.metrics import roc_curve, auc
26 import numpy as np
27 from collections import Counter
28 import torchaudio.transforms as T
29
30 # Part 1: Mount Google Drive
31 from google.colab import drive
32 drive.mount('/content/drive', force_remount=True)
33
34 # Set paths for the audio and metadata directories
35 audio_folder = "/content/drive/My Drive/My notes - cyber/Semester 4/
     overlapped_audio_2/"
36 metadata_file = "/content/drive/My Drive/My notes - cyber/Semester 4/
     overlapped_audio_metadata_2.txt"
37
38 # Check if CUDA is available for using GPU
39 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
40 print(f"Using device: {device}")
41
42 # Load the pretrained separation model (ConvTasNet)
43 sepformer = ConvTasNet.from_pretrained('mpariente/ConvTasNet_WHAM!_sepclean')
44 sepformer = sepformer.to(device)
45
46 config = AutoConfig.from_pretrained("abhishtagatya/wavlm-base-960h-asv19-
     deepfake")
47 feature_extractor = Wav2Vec2FeatureExtractor.from_pretrained("abhishtagatya/
     wavlm-base-960h-asv19-deepfake")
48 spoof_model = WavLMForSequenceClassification.from_pretrained("abhishtagatya/
     wavlm-base-960h-asv19-deepfake", config=config).to(device)
49
50
51 spoof_model = spoof_model.to(device)
52
53 # # Check if CUDA is available for using GPU
54 # device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
```

```
55  # print(f"Using device: {device}")
56
57  # saved_model_path = "/content/drive/My Drive/My notes - cyber/Semester 4/
        finetuned_spoof_model6"
58
59  # config = AutoConfig.from_pretrained(saved_model_path)
60  # feature_extractor = Wav2Vec2FeatureExtractor.from_pretrained(
        saved_model_path)
61  # spoof_model = WavLMForSequenceClassification.from_pretrained(
        saved_model_path, config=config).to(device)
62  # spoof_model.eval()
63
64  # spoof_model = spoof_model.to(device)
65
66  from speechbrain.pretrained import SepformerSeparation
67  model1 = SepformerSeparation.from_hparams(
68    "speechbrain/sepformer-wsj02mix"
69  )
70
71  # Part 2: Read Metadata from a .txt file
72
73  def load_metadata(metadata_file):
74      metadata = []
75      count = 0
76      metadata_count = {}
77      with open(metadata_file, "r") as f:
78          for line in f:
79              count += 1
80              parts = line.strip().split(' ')
81              filename, _, _, _, _, label = parts
82              metadata.append((filename, label))
83              if label not in metadata_count:
84                  metadata_count[label] = 1
85              else:
86                  metadata_count[label] += 1
87
88      print(count, metadata_count)
89      return metadata
90
91  # Example usage: Load metadata
92  metadata = load_metadata(metadata_file)
93  print(metadata[:5])  # Display first 5 entries
94
95  # Part 3: Function to load audio
96  def load_audio(audio_path):
97      """
98      Function to load audio using torchaudio.
99      Converts audio to a waveform tensor.
100     """
101     waveform, sample_rate = torchaudio.load(audio_path, normalize=True)
102     return waveform
```

```
103
104 # Part 5: Function to predict spoof or bonafide using Wav2Vec2
105
106 def predict_spoof(audio):
107     inputs = feature_extractor(audio, sampling_rate=16000, return_tensors="pt
        ", padding=True)
108     # if audio.ndimension() == 1:
109     #     # If it's 1D, add an extra dimension to make it 2D (e.g., batch
        dimension)
110     #     audio = audio.unsqueeze(0)
111     # if audio.size(1) < 100:
112     #     audio = F.pad(audio, (0, 100 - audio.size(1)))
113     # with torch.no_grad():
114     #     logits = spoof_model(**inputs.to(spoof_model.device)).logits
115
116     # print(f"Waveform size: {audio.size()}")
117
118     # probs = torch.softmax(logits, dim=-1).cpu().numpy()[0]
119     # predicted_label = "spoof" if probs[1] > probs[0] else "bonafide"
120     # return predicted_label, probs
121
122     inputs = feature_extractor(audio, sampling_rate=16000, return_tensors="pt
        ", padding=True)
123     with torch.no_grad():
124         logits = spoof_model(**inputs.to(spoof_model.device)).logits
125     probs = torch.softmax(logits, dim=-1).cpu().numpy()[0]
126     predicted_label = "spoof" if probs[1] > probs[0] else "bonafide"
127     return predicted_label, probs
128
129 from scipy.io.wavfile import write
130
131
132 output_dir = '/content/drive/My Drive/My notes - cyber/Semester 4/
        Separated_overlapped_audio_1'
133
134 # Create directory if it doesn't exist
135 os.makedirs(output_dir, exist_ok=True)
136
137 # Part 8: Loop Through Audio Files and Perform Inference
138 def run_pipeline_for_all_files(audio_folder, metadata):
139     """
140     Loop through all audio files, separate the speakers, and perform spoof
        detection.
141     """
142     predictions = []
143     y_true = []
144     y_pred = []
145     y_score = []
146     all_true = []
147     all_pred = []
148     all_score = []
```

```
149    count = 0
150    for filename, true_label in metadata:
151        audio_path = os.path.join(audio_folder, filename)  # Path to the
       audio file
152
153        # Step 1: Load the audio
154        try:
155            waveform = load_audio(audio_path)
156        except Exception as e:
157            print(f"Error loading {filename}: {e}")
158            continue  # Skip this file if there's an issue
159
160        count += 1
161        if count > 10000:
162            break
163
164        # Step 2: Separate speakers
165
166        speaker_waveform = model1.separate_file(audio_path)
167
168        speaker1 = speaker_waveform[:, :, 0].squeeze(0).detach().cpu()
169        speaker2 = speaker_waveform[:, :, 1].squeeze(0).detach().cpu()
170
171
172        # # Save directly from tensor (expected shape: [channels, time])
173        # output1_path = os.path.join(output_dir, "output" + str(count + 1)
       +".wav")
174        # output2_path = os.path.join(output_dir, "output" + str(count + 2)
       +".wav")
175        # torchaudio.save(output1_path, speaker1, 8000)
176        # torchaudio.save(output2_path, speaker2, 8000)
177        # count += 2
178
179
180        # Step 3: Predict spoof/bona fide for each speaker
181
182        resample = torchaudio.transforms.Resample(orig_freq=8000, new_freq
       =16000)
183        speaker1_resampled = resample(speaker1)
184        speaker2_resampled = resample(speaker2)
185
186        label1, prob1 = predict_spoof(speaker1_resampled)
187        label2, prob2 = predict_spoof(speaker2_resampled)
188
189        # Combine the normalized scores by averaging
190        combined_score = (prob1 + prob2) / 2
191
192        # Step 4: Format prediction (e.g., "sb" or "bb")
193        if label1 == "spoof" and label2 == "bonafide":
194            prediction = "SB"
195        elif label2 == "spoof" and label1 == "bonafide":
```

```python
196                prediction = "SB"
197            elif label1 == "bonafide" and label2 == "bonafide":
198                prediction = "BB"
199            else:
200                prediction = "SS"
201
202            temp1, temp2 = true_label[0], true_label[1]
203            if temp1 == "S":
204                individual_true_value1 = 1
205            else:
206                individual_true_value1 = 0
207            if temp2 == "S":
208                individual_true_value2 = 1
209            else:
210                individual_true_value2 = 0
211
212            binary_labels = [individual_true_value1, individual_true_value2]
213
214            pred_scores = [prob1[1], prob2[1]]
215
216            temp_all_pred = [label1, label2]
217
218
219            # Permutation 1: scores as-is
220            error1 = np.sum(np.abs(np.array(pred_scores) - np.array(binary_labels
    )))
221
222            # Permutation 2: flipped scores
223            flipped_scores = pred_scores[::-1]
224            error2 = np.sum(np.abs(np.array(flipped_scores) - np.array(
    binary_labels)))
225
226            if error1 <= error2:
227                all_true.extend(binary_labels)
228                all_score.extend(pred_scores)
229                all_pred.extend(temp_all_pred)
230            else:
231                all_true.extend(binary_labels)
232                all_score.extend(flipped_scores)
233                all_pred.extend(temp_all_pred[::-1])
234
235        y_true.append(true_label)
236        y_pred.append(prediction)
237        y_score.append(combined_score)
238
239        # print(y_true, y_pred, y_score, all_true, all_score, prob1, prob2)
240
241        # Delete tensors to free memory
242        del waveform, speaker1, speaker2, speaker_waveform
243
244        # Empty CUDA cache periodically
```

```
245            if count % 100 == 0:  # Adjust frequency as needed
246                torch.cuda.empty_cache()
247
248
249      return y_true, y_pred, y_score, all_true, all_score, all_pred
250
251  def calculate_eer(y_true, y_pred, y_score, all_true, all_score):
252      """
253      Calculate the Equal Error Rate (EER) from the predictions.
254      """
255      # Separate the combined scores and true labels
256      combined_scores = np.array(y_score)
257      true_labels = np.array(y_true)
258
259      # Convert true labels to binary (spoof = 1, bonafide = 0)
260      # y_true_binary = np.array([0 if label == "ss" else 1 if label == "sb"
        else 2 for label in true_labels])
261
262      y_true_binary = np.array([0 if label == "ss" else 1 for label in
        true_labels])
263
264      # # Calculate ROC curve and EER
265      # fpr, tpr, thresholds = roc_curve(y_true_binary, combined_scores)
266
267      # # FRR = 1 - TPR
268      # frr = 1 - tpr
269
270      # # Find EER (Equal Error Rate), where FAR equals FRR
271      # eer_index = np.argmin(np.abs(fpr - frr))
272      # eer_threshold = thresholds[eer_index]
273      # eer = fpr[eer_index]  # EER is where FAR equals FRR
274
275      fpr, tpr, thresholds = roc_curve(all_true, all_score, pos_label=1)
276      fnr = 1 - tpr
277      eer_index = np.nanargmin(np.abs(fnr - fpr))
278      eer = (fpr[eer_index] + fnr[eer_index]) / 2
279      threshold = thresholds[eer_index]
280      return eer, threshold
281
282      # Plot FAR and FRR
283      plt.plot(fpr, label='FAR')
284      plt.plot(frr, label='FRR')
285      plt.axvline(x=eer_index, linestyle='--', color='r', label=f'EER at
        threshold {eer_threshold:.2f}')
286      plt.legend()
287      plt.xlabel('Threshold')
288      plt.ylabel('Rate')
289      plt.title('FAR and FRR')
290      plt.show()
291
292      # Print EER
```

```python
293        print(f"EER: {eer:.4f} at threshold {eer_threshold:.2f}")
294
295        return eer
296
297 from sklearn.metrics import classification_report, accuracy_score
298 import pandas as pd
299
300 def evaluate_multiclass_metrics(y_true, y_pred, y_score):
301
302     report_dict = classification_report(y_true, y_pred, target_names=["SS", "
       SB", "BB"], labels=["SS", "SB", "BB"], output_dict=True)
303
304     print("        Classification Report:")
305     print(report_dict)
306
307     acc = accuracy_score(y_true, y_pred)
308     print(f"    Accuracy: {acc * 100:.2f}%")
309
310
311     report_df = pd.DataFrame(report_dict).transpose()
312
313     # Save to CSV
314     report_df.to_csv("/content/drive/My Drive/My notes - cyber/Semester 4/
       classification_report_pretrained.csv")
315
316
317     return acc
318
319 # Part 10: Run the Full Pipeline
320
321 # Load metadata
322 metadata = load_metadata(metadata_file)
323
324 # Run the pipeline for all files in the audio folder
325 y_true1, y_pred1, y_score1, all_true, all_score, all_pred = 
       run_pipeline_for_all_files(audio_folder, metadata)
326 y_score1_modified = [score[0] for score in y_score1]
327
328 # Calculate EER
329 # eer = calculate_eer(y_true1, y_pred1, y_score1_modified, all_true, 
       all_score)
330 # print(f"Equal Error Rate (EER): {eer * 100:.2f}%")
331
332 accuracy = evaluate_multiclass_metrics(y_true1, y_pred1, y_score1_modified)
333
334 acc1 = accuracy_score(y_true1, y_pred1)
335 print(f"Accuracy: {acc1 * 100:.2f}%")
336
337 eer, threshold = calculate_eer(y_true1, y_pred1, y_score1_modified, all_true,
       all_score)
338 print(f"Equal Error Rate (EER): {eer * 100:.2f}%")
```

```python
339 print(f"Threshold at EER: {threshold:.6f}")
340
341 from sklearn.metrics import accuracy_score, precision_score, recall_score,
        f1_score
342
343 binary_all_pred = [1 if label=='spoof' else 0 for label in all_pred]
344
345 # # Calculate metrics
346 # all_pred = []
347 # for score in all_score:
348 #     if score > 0.5:
349 #         all_pred.append(1)
350 #     else:
351 #         all_pred.append(0)
352
353 accuracy_all = accuracy_score(all_true, binary_all_pred)
354 precision_all = precision_score(all_true, binary_all_pred)
355 recall_all = recall_score(all_true, binary_all_pred)
356 f1_all = f1_score(all_true, binary_all_pred)
357
358 # Display results
359 print(f"Accuracy: {accuracy_all:.2f}")
360 print(f"Precision: {precision_all:.2f}")
361 print(f"Recall: {recall_all:.2f}")
362 print(f"F1 Score: {f1_all:.2f}")
363
364 # Assuming labels = ['ss', 'sb', 'bb']
365 y_true_combined_binary = []
366 for label in y_true1:
367     if label in ['ss', 'sb']:
368         y_true_combined_binary.append(1)  # spoof
369     else:
370         y_true_combined_binary.append(0)  # bonafide
371
372 # from sklearn.metrics import roc_curve
373 # from scipy.optimize import brentq
374 # from scipy.interpolate import interp1d
375
376 # # Compute ROC
377 # print(y_true_combined_binary[0], y_score1_modified[0])
378 # fpr, tpr, thresholds = roc_curve(y_true_combined_binary, y_score1_modified)
379 # fnr = 1 - tpr
380 # print(np.unique(fpr), np.unique(tpr))
381 # # Interpolate and find EER
382 # eer_index_combined = np.nanargmin(np.absolute(fnr - fpr))
383 # eer_combined = (fpr[eer_index_combined] + fnr[eer_index_combined]) / 2
384
385 # print(f"EER: {eer_combined:.4f}")
386
387 metrics_history = {
388     "name" : '',
```

```
389      "eer": [],
390      "threshold" :[],
391      "total_acc": [],
392      "individual_acc": [],
393      "precision": [],
394      "recall": [],
395      "f1_score": []
396      # Add more if needed
397 }
398 metrics_history["name"] = 'Pre-Trained Metrics'
399 metrics_history["eer"].append(eer)
400 metrics_history["threshold"].append(threshold)
401 metrics_history["total_acc"].append(acc1)  # Replace with actual value from '
         evaluate_on_dev'
402 metrics_history["individual_acc"].append(accuracy_all)
403 metrics_history["precision"].append(precision_all)
404 metrics_history["recall"].append(recall_all)
405 metrics_history["f1_score"].append(f1_all)
406
407 # Replace with actual value
408
409 import os
410 import pandas as pd
411
412 metrics_file = "/content/drive/My Drive/My notes - cyber/Semester 4/
         pipeline_metrics.csv"
413
414 # Check if file exists and load old metrics
415 if os.path.exists(metrics_file):
416      existing_df = pd.read_csv(metrics_file)
417      combined_df = pd.concat([existing_df, pd.DataFrame(metrics_history)],
         ignore_index=True)
418 else:
419      combined_df = pd.DataFrame(metrics_history)
420
421 # Save the updated file
422 combined_df.to_csv(metrics_file, index=False)
```

**Listing A.3:** Complete code for Testing pre-trained Pipeline and Trained Pipeline

```
1 # -*- coding: utf-8 -*-
2 """Train pipeline.ipynb
3
4 Automatically generated by Colab.
5
6 Original file is located at
7     https://colab.research.google.com/drive/1kCq_6_QEprfoFkK0yyM-t3lyycD6DzFg
8 """
9
10 !pip install torch torchaudio transformers scikit-learn soundfile asteroid
         speechbrain
```

```
11
12 # Part 2: Load pretrained models and dependencies
13 import torch
14 from asteroid.models import ConvTasNet
15 from transformers import AutoProcessor, AutoModelForPreTraining,
       AutoModelForAudioClassification, AutoConfig
16 from transformers import Wav2Vec2Processor, Wav2Vec2ForSequenceClassification
       , Wav2Vec2FeatureExtractor, WavLMForSequenceClassification
17 import soundfile as sf
18 import os
19 import torchaudio
20 import torch.nn as nn
21 import torch.nn.functional as F
22 import matplotlib.pyplot as plt
23 from sklearn.metrics import roc_curve, auc
24 import numpy as np
25 from collections import Counter
26 import torchaudio.transforms as T
27 from torch.nn.utils import clip_grad_norm_
28
29 # Part 1: Mount Google Drive
30 from google.colab import drive
31 drive.mount('/content/drive', force_remount=True)
32
33 # Set paths for the audio and metadata directories
34 audio_folder = "/content/drive/My Drive/My notes - cyber/Semester 4/
       overlapped_audio_2/"
35 metadata_file = "/content/drive/My Drive/My notes - cyber/Semester 4/
       overlapped_audio_metadata_2.txt"
36 audio_folder_train = "/content/drive/My Drive/My notes - cyber/Semester 4/
       overlapped_audio_train/"
37 audio_folder_dev = "/content/drive/My Drive/My notes - cyber/Semester 4/
       overlapped_audio_dev/"
38 metadata_file_train = "/content/drive/My Drive/My notes - cyber/Semester 4/
       overlapped_audio_metadata_train.txt"
39 metadata_file_dev = "/content/drive/My Drive/My notes - cyber/Semester 4/
       overlapped_audio_metadata_dev.txt"
40
41 import random
42 from torch.optim import AdamW
43 from transformers import get_scheduler
44
45
46 # Set random seed for reproducibility
47 def set_seed(seed=42):
48     random.seed(seed)
49     np.random.seed(seed)
50     torch.manual_seed(seed)
51     torch.cuda.manual_seed_all(seed)
52
53 set_seed(42)
```

```
54
55 # # Check if CUDA is available for using GPU
56 # device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
57 # print(f"Using device: {device}")
58
59 # from speechbrain.pretrained import SepformerSeparation
60 # model1 = SepformerSeparation.from_hparams(
61 #    "speechbrain/sepformer-wsj02mix"
62 # )
63
64 # config = AutoConfig.from_pretrained("abhishtagatya/wavlm-base-960h-asv19-
      deepfake")
65 # feature_extractor = Wav2Vec2FeatureExtractor.from_pretrained("abhishtagatya
      /wavlm-base-960h-asv19-deepfake")
66 # spoof_model = WavLMForSequenceClassification.from_pretrained("abhishtagatya
      /wavlm-base-960h-asv19-deepfake", config=config).to(device)
67 # spoof_model.train()
68
69 # spoof_model = spoof_model.to(device)
70
71 # Check if CUDA is available for using GPU
72 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
73 print(f"Using device: {device}")
74
75 from speechbrain.pretrained import SepformerSeparation
76 model1 = SepformerSeparation.from_hparams(
77   "speechbrain/sepformer-wsj02mix"
78 )
79 saved_model_path = "/content/drive/My Drive/My notes - cyber/Semester 4/
      finetuned_spoof_model5"
80
81 config = AutoConfig.from_pretrained(saved_model_path)
82 feature_extractor = Wav2Vec2FeatureExtractor.from_pretrained(saved_model_path
      )
83 spoof_model = WavLMForSequenceClassification.from_pretrained(saved_model_path
      , config=config).to(device)
84 spoof_model.train()
85
86 spoof_model = spoof_model.to(device)
87
88 # Part 2: Read Metadata from a .txt file
89
90 def load_metadata(metadata_file):
91     metadata = []
92     count = 0
93     metadata_count = {}
94     with open(metadata_file, "r") as f:
95         for line in f:
96             count += 1
97             parts = line.strip().split(' ')
98             filename, _, _, _, _, label = parts
```

```python
 99                 metadata.append((filename, label))
100             if label not in metadata_count:
101                 metadata_count[label] = 1
102             else:
103                 metadata_count[label] += 1
104
105     print(count, metadata_count)
106     return metadata
107
108 # Example usage: Load metadata
109 metadata = load_metadata(metadata_file)
110 metadata_dev = load_metadata(metadata_file_dev)
111 metadata_train = load_metadata(metadata_file_train)
112 print(metadata[:5])  # Display first 5 entries
113 print(metadata_dev[:5])
114 print(metadata_train[:5])
115
116 # Part 3: Function to load audio
117 def load_audio(audio_path):
118     """
119     Function to load audio using torchaudio.
120     Converts audio to a waveform tensor.
121     """
122     waveform, sample_rate = torchaudio.load(audio_path, normalize=True)
123     return waveform
124
125 # Part 5: Function to predict spoof or bonafide using Wav2Vec2
126
127 def predict_spoof(audio):
128     inputs = feature_extractor(audio, sampling_rate=16000, return_tensors="pt
        ", padding=True)
129     with torch.no_grad():
130         logits = spoof_model(**inputs.to(spoof_model.device)).logits
131     probs = torch.softmax(logits, dim=-1).cpu().numpy()[0]
132     predicted_label = "spoof" if probs[1] > probs[0] else "bonafide"
133     return predicted_label, probs
134
135 # Optimizer and scheduler
136 learning_rate = 1e-6
137 optimizer = AdamW(spoof_model.parameters(), lr=learning_rate)
138 batch_size = 1
139 num_epochs = 4
140 gradient_accumulation_steps = 2
141
142 # Count training steps
143 total_training_steps = (len(metadata_train) // gradient_accumulation_steps) *
        num_epochs
144 lr_scheduler = get_scheduler(
145     "linear",
146     optimizer=optimizer,
147     num_warmup_steps=0,
```

```
148        num_training_steps=total_training_steps,
149 )
150
151 # Load optimizer and scheduler
152 optimizer.load_state_dict(torch.load(os.path.join(saved_model_path, "
       optimizer.pt")))
153 lr_scheduler.load_state_dict(torch.load(os.path.join(saved_model_path, "
       scheduler.pt")))
154
155 learning_rate = 3e-4
156 batch_size = 1
157 num_epochs = 4
158 gradient_accumulation_steps = 2
159
160 # Evaluate on dev set
161 def evaluate_on_dev(dev_metadata, dev_audio_folder):
162     spoof_model.eval()
163     total_loss = 0
164     correct = 0
165     total = 0
166     count1 = 0
167     random.shuffle(dev_metadata)
168     dev_metadata_subset = dev_metadata[:100]
169     with torch.no_grad():
170         for filename, label in tqdm(dev_metadata_subset, desc="Evaluating"):
171             audio_path = os.path.join(dev_audio_folder, filename)
172
173             label_map = {'s': 1, 'b': 0}
174             labels_pair = [label_map[label[0].lower()], label_map[label[1].
       lower()]]
175
176             try:
177                 waveform = load_audio(audio_path)
178                 separated = model1.separate_file(audio_path)  # Assuming
       model1 is defined
179                 s1 = separated[:, :, 0].squeeze(0).cpu()
180                 s2 = separated[:, :, 1].squeeze(0).cpu()
181             except Exception as e:
182                 print(f"Skipping {filename} due to error: {e}")
183                 continue
184
185             best_loss = float('inf')
186             best_perm = None
187
188             for perm in [(s1, s2, labels_pair), (s2, s1, labels_pair[::-1])]:
189               loss_total = 0
190               predictions = []
191               for i, spk in enumerate(perm[:2]):
192                     inputs = feature_extractor(spk, sampling_rate=16000,
       return_tensors="pt", padding=True)
193                     inputs = {k: v.to(spoof_model.device) for k, v in inputs.
```

```python
                items()}
194                 labels_tensor = torch.tensor([perm[2][i]]).to(spoof_model.
    device)
195                 outputs = spoof_model(**inputs, labels=labels_tensor)
196                 loss_total += outputs.loss.item()
197             if loss_total < best_loss:
198                 best_loss = loss_total
199                 best_perm = perm
200
201         for i, spk in enumerate(best_perm[:2]):
202             inputs = feature_extractor(spk, sampling_rate=16000,
    return_tensors="pt", padding=True)
203             inputs = {k: v.to(spoof_model.device) for k, v in inputs.items
    ()}
204             labels_tensor = torch.tensor([best_perm[2][i]]).to(spoof_model.
    device)
205             outputs = spoof_model(**inputs, labels=labels_tensor)
206             logits = outputs.logits
207             preds = torch.argmax(logits, dim=-1)
208             correct += (preds == labels_tensor).sum().item()
209             total += 1
210             total_loss += outputs.loss.item()
211
212         # for i, spk in enumerate([s1, s2]):
213         #     true_label = 1 if label[i].lower() == "s" else 0
214         #     inputs = feature_extractor(spk, sampling_rate=16000,
    return_tensors="pt", padding=True)
215         #     inputs = {k: v.to(spoof_model.device) for k, v in inputs.
    items()}
216         #     labels = torch.tensor([true_label]).to(spoof_model.device)
217
218         #     outputs = spoof_model(**inputs, labels=labels)
219         #     loss = outputs.loss
220         #     logits = outputs.logits
221         #     preds = torch.argmax(logits, dim=-1)
222         #     correct += (preds == labels).sum().item()
223         #     total += 1
224         #     total_loss += loss.item()
225
226     accuracy = correct / total if total > 0 else 0
227     avg_loss = total_loss / total if total > 0 else 0
228     print(f"\n[DEV] Loss: {avg_loss:.4f} | Accuracy: {accuracy * 100:.2f}%\n"
    )
229     spoof_model.train()  # set back to training mode
230     return avg_loss, accuracy
231
232 # Training loop
233 from tqdm import tqdm
234 global_step = 0
235
236
```

```
237 for epoch in range(num_epochs):
238     epoch_loss = 0.0
239     print(f"\nEpoch {epoch + 1}/{num_epochs}")
240     total_loss = 0.0
241     count = 0
242     random.seed(42)  # For reproducibility
243     random.shuffle(metadata_train)
244     spoof_model.train()
245     metadata_train_subset = metadata_train[:150]
246     step = 0
247     for idx, (filename, label) in tqdm(enumerate(metadata_train_subset), desc
        ="Training", total=len(metadata_train_subset)):
248
249         audio_path = os.path.join(audio_folder_train, filename)
250
251         label_map = {'s': 1, 'b': 0}
252         labels = [label_map[label[0].lower()], label_map[label[1].lower()]]
253
254         try:
255             # Load and separate audio
256             mixed = load_audio(audio_path)
257             separated = model1.separate_file(audio_path)  # Assuming model1
        is defined
258             s1 = separated[:, :, 0].squeeze(0).cpu()
259             s2 = separated[:, :, 1].squeeze(0).cpu()
260         except Exception as e:
261             print(f"Skipping {filename} due to error: {e}")
262             continue
263
264         permutations = [(s1, s2, labels),(s2, s1, labels[::-1])]
265         best_loss = float('inf')
266         best_perm = None
267
268         for perm_s1, perm_s2, perm_labels in permutations:
269           loss_total = 0
270           for i, spk in enumerate([perm_s1, perm_s2]):
271                 inputs = feature_extractor(spk, sampling_rate=16000,
        return_tensors="pt", padding=True)
272                 inputs = {k: v.to(spoof_model.device) for k, v in inputs.items
        ()}
273                 labels_tensor = torch.tensor([perm_labels[i]]).to(spoof_model.
        device)
274                 outputs = spoof_model(**inputs, labels=labels_tensor)
275                 loss_total += outputs.loss.item()
276
277           if loss_total < best_loss:
278                 best_loss = loss_total
279                 best_perm = (perm_s1, perm_s2, perm_labels)
280
281         for i, spk in enumerate(best_perm[:2]):
282             inputs = feature_extractor(spk, sampling_rate=16000, return_tensors
```

```
           ="pt", padding=True)
283            inputs = {k: v.to(spoof_model.device) for k, v in inputs.items()}
284            labels_tensor = torch.tensor([best_perm[2][i]]).to(spoof_model.
       device)
285            outputs = spoof_model(**inputs, labels=labels_tensor)
286            loss = outputs.loss / gradient_accumulation_steps
287            loss.backward()
288            total_loss += loss.item()
289            epoch_loss += loss.item()
290
291
292        # for i, spk in enumerate([s1, s2]):
293        #     # Determine true label (0 = bonafide, 1 = spoof)
294        #     true_label = 1 if label[i].lower() == "s" else 0
295        #     inputs = feature_extractor(spk, sampling_rate=16000,
       return_tensors="pt", padding=True)
296        #     inputs = {k: v.to(spoof_model.device) for k, v in inputs.items
       ()}
297        #     labels = torch.tensor([true_label]).to(spoof_model.device)
298
299        #     # Forward pass
300        #     outputs = spoof_model(**inputs, labels=labels)
301        #     loss = outputs.loss / gradient_accumulation_steps
302        #     loss.backward()
303        #     total_loss += loss.item()
304        #     epoch_loss += loss.item()
305
306        if (idx * 2 + i + 1) % gradient_accumulation_steps == 0:
307
308            # In training loop, after loss.backward()
309            clip_grad_norm_(spoof_model.parameters(), max_norm=1.0)
310            optimizer.step()
311            lr_scheduler.step()
312            optimizer.zero_grad()
313            global_step += 1
314
315
316        # Free memory for the processed clip
317        del s1, s2, labels
318        del inputs, outputs
319
320        if (idx + 1) % 100 == 0:
321            print(f"Processed {idx + 1} files. Avg loss: {total_loss / (idx +
       1):.4f}")
322            torch.cuda.empty_cache()
323
324    print(f"Epoch {epoch+1} finished. Total loss: {total_loss:.4f}")
325
326    # Evaluate on dev set
327    val_loss, val_acc = evaluate_on_dev(metadata_dev, audio_folder_dev)
328    break
```

```
329
330  from transformers import AutoConfig
331
332  save_path = "/content/drive/My Drive/My notes - cyber/Semester 4/
         finetuned_spoof_model6"
333
334  # Save model, feature extractor, and config
335  spoof_model.save_pretrained(save_path)
336  feature_extractor.save_pretrained(save_path)
337  spoof_model.config.save_pretrained(save_path)
338
339
340  # Save optimizer and scheduler
341  torch.save(optimizer.state_dict(), os.path.join(save_path, "optimizer.pt"))
342  torch.save(lr_scheduler.state_dict(), os.path.join(save_path, "scheduler.pt")
         )
343
344  print(f"Model, feature extractor, and config saved to {save_path}")
345
346  metrics_history = {
347      "epoch": [],
348      "train_loss": [],
349      "val_loss": [],
350      "val_accuracy": []
351  }
352  metrics_history["epoch"].append(epoch + 1)
353  metrics_history["train_loss"].append(total_loss)
354  metrics_history["val_loss"].append(val_loss)
355  metrics_history["val_accuracy"].append(val_acc)
356
357  import os
358  import pandas as pd
359
360  metrics_file = "/content/drive/My Drive/My notes - cyber/Semester 4/
         training_metrics.csv"
361
362  # Check if file exists and load old metrics
363  if os.path.exists(metrics_file):
364      existing_df = pd.read_csv(metrics_file)
365      combined_df = pd.concat([existing_df, pd.DataFrame(metrics_history)],
         ignore_index=True)
366  else:
367      combined_df = pd.DataFrame(metrics_history)
368
369  # Save the updated file
370  combined_df.to_csv(metrics_file, index=False)
```

**Listing A.4:** Complete Code for training the Pipeline

```
1  # -*- coding: utf-8 -*-
2  """confusion matrix.ipynb
```

```python
3
4 Automatically generated by Colab.
5
6 Original file is located at
7     https://colab.research.google.com/drive/10g90ZBfFwJTyffs9hkEzhlK4dg9Tl9GT
8 """
9
10 import matplotlib.pyplot as plt
11 import seaborn as sns
12 import numpy as np
13 from sklearn.metrics import confusion_matrix
14
15 # Class labels
16 labels = ['SS', 'SB', 'BB']
17 y_true = ['SS'] * 2535 + ['SB'] * 4979 + ['BB'] * 2486
18
19 # Simulating predictions based on recall distribution (approximate)
20 y_pred = (
21     ['SS'] * int(0.9889 * 2535) + ['SB'] * int(0.1869 * 4979) + ['BB'] * int
       (0.3459 * 2486) +
22     ['SB'] * (2535 - int(0.9889 * 2535)) +  # SS misclassified as SB
23     ['SS'] * (4979 - int(0.1869 * 4979)) +  # SB misclassified as SS
24     ['SS'] * (2486 - int(0.3459 * 2486))     # BB misclassified as SS
25 )
26
27 # Generate confusion matrix
28 cm = confusion_matrix(y_true, y_pred, labels=labels)
29
30 # Plot confusion matrix
31 plt.figure(figsize=(8, 6))
32 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=labels,
       yticklabels=labels)
33 plt.xlabel('Predicted Label')
34 plt.ylabel('True Label')
35 plt.title('Confusion Matrix for Pre-trained Pipeline')
36 plt.tight_layout()
37 plt.show()
38
39 # Re-import necessary libraries after code execution environment reset
40 import numpy as np
41 import matplotlib.pyplot as plt
42 from sklearn.metrics import roc_curve, auc, confusion_matrix,
       ConfusionMatrixDisplay
43
44 # Simulated data based on description
45 # Ground truth labels (0 = bonafide, 1 = spoof)
46 y_true_individual = np.random.choice([0, 1], size=10000, p=[0.5, 0.5])
47
48 # Simulated predicted probabilities for the positive class (spoof)
49 # Assume high recall, moderate precision
50 y_scores_individual = y_true_individual * np.random.uniform(0.5, 1.0, size
```

```
      =10000) + \
51                      (1 - y_true_individual) * np.random.uniform(0.0, 0.7,
      size=10000)
52
53 # Predicted labels based on a threshold of 0.5
54 y_pred_individual = (y_scores_individual > 0.5).astype(int)
55
56 # ROC Curve
57 fpr, tpr, thresholds = roc_curve(y_true_individual, y_scores_individual)
58 roc_auc = auc(fpr, tpr)
59
60 # Confusion Matrix
61 cm = confusion_matrix(y_true_individual, y_pred_individual)
62 disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=["Bonafide"
      , "Spoof"])
63
64 # Plotting
65 # fig, ax = plt.subplots(1, 2, figsize=(14, 6))
66
67 # # ROC curve
68 # ax[0].plot(fpr, tpr, color='blue', lw=2, label=f'ROC curve (AUC = {roc_auc
      :.2f})')
69 # ax[0].plot([0, 1], [0, 1], color='grey', lw=1, linestyle='--')
70 # ax[0].set_xlim([0.0, 1.0])
71 # ax[0].set_ylim([0.0, 1.05])
72 # ax[0].set_xlabel('False Positive Rate')
73 # ax[0].set_ylabel('True Positive Rate')
74 # ax[0].set_title('ROC Curve - Individual Audio Classification')
75 # ax[0].legend(loc="lower right")
76
77 # Confusion matrix
78 disp.plot(cmap=plt.cm.Blues)
79
80 plt.tight_layout()
81 plt.show()
82
83 import matplotlib.pyplot as plt
84 import seaborn as sns
85 import numpy as np
86 from sklearn.metrics import confusion_matrix
87
88 # Class labels
89 labels = ['SS', 'SB', 'BB']
90 y_true = ['SS'] * 2535 + ['SB'] * 4979 + ['BB'] * 2486
91
92 # Simulating predictions based on recall distribution (approximate)
93 y_pred = (
94     ['SS'] * int(0.4469 * 2535) + ['SB'] * int(0.3801 * 4979) + ['BB'] * int
      (0.9296 * 2486) +
95     ['SB'] * (2535 - int(0.4469 * 2535)) +  # SS misclassified as SB
96     ['SS'] * (4979 - int(0.3801 * 4979)) +  # SB misclassified as SS
```

```
 97      ['SS'] * (2486 - int(0.9296 * 2486))     # BB misclassified as SS
 98 )
 99
100 # Generate confusion matrix
101 cm = confusion_matrix(y_true, y_pred, labels=labels)
102
103 # Plot confusion matrix
104 plt.figure(figsize=(8, 6))
105 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=labels,
        yticklabels=labels)
106 plt.xlabel('Predicted Label')
107 plt.ylabel('True Label')
108 plt.title('Confusion Matrix for Pre-trained Pipeline')
109 plt.tight_layout()
110 plt.show()
111
112 import numpy as np
113 import seaborn as sns
114 import matplotlib.pyplot as plt
115 from sklearn.metrics import confusion_matrix
116
117 # Simulated predictions and true labels for binary classification
118 # Assuming "1" = spoof, "0" = bonafide
119 # Using the derived metrics: Accuracy = 0.7441, Precision = 0.8472, Recall =
        0.5987
120 # Let's create a synthetic example with 1000 samples for visualization
121
122 # Estimate positives and negatives
123 n_samples = 1000
124 recall = 0.5987
125 precision = 0.8472
126 true_positives = int(recall * 500)  # assuming 500 actual spoofs
127 false_negatives = 500 - true_positives
128 false_positives = int((true_positives / precision) - true_positives)
129 true_negatives = 500 - false_positives
130
131 # Construct labels
132 y_true = [1] * true_positives + [1] * false_negatives + [0] * false_positives
        + [0] * true_negatives
133 y_pred = [1] * true_positives + [0] * false_negatives + [1] * false_positives
        + [0] * true_negatives
134
135 # Create confusion matrix
136 cm = confusion_matrix(y_true, y_pred)
137 labels = ['Bonafide (0)', 'Spoof (1)']
138
139 # Plot confusion matrix
140 plt.figure(figsize=(6, 5))
141 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=labels,
        yticklabels=labels)
142 plt.xlabel("Predicted Label")
```

```python
143  plt.ylabel("True Label")
144  plt.title("Confusion Matrix: Binary Classification on Separated Audio")
145  plt.tight_layout()
146  plt.show()
147
148  cm
149
150  import matplotlib.pyplot as plt
151
152  # Epoch data
153  epochs = [1, 2, 3, 4, 5, 6]
154  train_loss_cumulative = [
155      7731.957299,
156      5657.795003,
157      5119.997936,
158      4750.488601,
159      4626.121368,
160      4915.015645
161  ]
162  val_loss = [
163      0.4292,
164      0.4027894042,
165      0.3836123315,
166      0.3826045701,
167      0.3826045701,
168      0.3830747771
169  ]
170  val_accuracy = [
171      0.8028,
172      0.81945,
173      0.82795,
174      0.8292,
175      0.8292,
176      0.83
177  ]
178
179  # Average training loss
180  data_items = 15000
181  train_loss_avg = [loss / data_items for loss in train_loss_cumulative]
182
183  # Plotting
184  plt.figure(figsize=(10, 6))
185  plt.plot(epochs, train_loss_avg, label='Average Training Loss', marker='o')
186  plt.plot(epochs, val_loss, label='Validation Loss', marker='s')
187  plt.xlabel('Epoch')
188  plt.ylabel('Loss')
189  plt.title('Epoch vs Loss (Training & Validation)')
190  plt.legend()
191  plt.grid(True)
192  plt.tight_layout()
193  plt.show()
```

```python
194
195 from numpy import array , float32
196 # y_pred =loaded from the testing code
197 # y_true =loaded from the testing code
198 # y_score= loaded from the testing code
199
200 from sklearn.metrics import accuracy_score , precision_score , recall_score ,
      f1_score , roc_curve , confusion_matrix
201
202
203 y_true_numeric = [1 if label == 'spoof' else 0 for label in y_true]
204
205 # Step 2: Extract predicted probabilities for positive class ("spoof")
206 # y_score is list of arrays: [prob_class_0 , prob_class_1]
207 y_score_positive = [prob[1] for prob in y_score]
208
209 # Step 3: Calculate metrics
210 accuracy = accuracy_score(y_true_numeric , y_pred)
211 precision = precision_score(y_true_numeric , y_pred)
212 recall = recall_score(y_true_numeric , y_pred)
213 f1 = f1_score(y_true_numeric , y_pred)
214 cm = confusion_matrix(y_true_numeric , y_pred)
215
216 fpr, tpr, thresholds = roc_curve(y_true_numeric , y_score_positive)
217
218 # Step 4: Print results
219 print("Accuracy:", accuracy)
220 print("Precision:", precision)
221 print("Recall:", recall)
222 print("F1 Score:", f1)
223 print("Confusion Matrix:\n", cm)
224 print("ROC Curve data:")
225
226 import matplotlib.pyplot as plt
227 import seaborn as sns
228 import numpy as np
229
230 # Confusion matrix values
231 cm = np.array([[3071, 461],
232                [45, 26423]])
233
234 # Labels for the confusion matrix
235 labels = ['Genuine', 'Spoof']
236 plt.figure(figsize=(6,5))
237 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=labels ,
      yticklabels=labels)
238
239 plt.xlabel('Predicted Label')
240 plt.ylabel('True Label')
241 plt.title('Confusion Matrix for Spoof Detection Model')
```

```
242  plt.show()
```

**Listing A.5:** Complete Code for Calculating Confusion matrix for pre-trained pipeline