# Momenta - Audio Deepfake Detection Take-Home Assessment

## Part 1: Research & Selection

After reviewing the GitHub repository on audio deepfake detection ,I have identified three promising approaches for detecting AI-generated human speech with potential for real-time analysis in conversational settings:

1. **AASIST: Audio Anti-Spoofing Using Integrated Spectro-Temporal Graph Attention Networks**

   ○ **Key Technical Innovation:** AASIST employs graph attention networks to model both spectral and temporal relationships in audio data, enhancing the detection of subtle artifacts introduced by deepfake generation methods.GitHub

   ○ **Reported Performance Metrics:** Achieved an Equal Error Rate (EER) of 0.83% and a minimum tandem Detection Cost Function (t-DCF) of 0.028 on the Logical Access (LA) scenario.GitHub

   ○ **Why Promising:** The model's ability to capture intricate spectro-temporal patterns makes it adept at identifying AI-generated speech, which often contains nuanced inconsistencies not present in genuine audio.

   ○ **Potential Limitations or Challenges:** The computational complexity of graph attention networks may pose challenges for real-time deployment, necessitating optimization strategies to maintain efficiency without compromising accuracy.

2. **End-to-End Dual-Branch Network for Synthetic Speech Detection**

   ○ **Key Technical Innovation:** This approach integrates Linear Frequency Cepstral Coefficients (LFCC) and Constant-Q Transform (CQT) features within a dual-branch network architecture, leveraging both time-domain and frequency-domain information for improved detection.GitHub

   ○ **Reported Performance Metrics:** Recorded an EER of 0.80% and a t-DCF of 0.021 in the LA scenario.GitHub+1Resemble AI+1

   ○ **Why Promising:** By combining complementary audio features, the model enhances its capability to discern synthetic speech patterns, which is crucial for accurate detection in real conversational contexts.GitHub

- 
  - 
    - **Potential Limitations or Challenges:** The dual-branch structure may increase model complexity and computational load, potentially impacting real-time processing capabilities unless optimized effectively.

  3. **RawNet2 with Sinc Filters for End-to-End Anti-Spoofing**

     - **Key Technical Innovation:** RawNet2 utilizes Sinc filters for efficient feature extraction directly from raw audio waveforms, facilitating end-to-end learning without the need for handcrafted features.GitHub

     - **Reported Performance Metrics:** Achieved an EER of 1.12% and a t-DCF of 0.033 in the LA scenario.GitHub+1Resemble AI+1

     - **Why Promising:** The end-to-end nature and use of Sinc filters allow for efficient processing, making it suitable for real-time applications where rapid detection of AI-generated speech is required.

     - **Potential Limitations or Challenges:** While efficient, the model's performance may vary with different types of synthetic speech, necessitating continuous updates and training with diverse datasets to maintain robustness.


I'll implement **RawNet2 with Sinc Filters for End-to-End Anti-Spoofing** because:

- It processes raw audio directly, reducing reliance on handcrafted features.

- It has a relatively lightweight architecture, making it suitable for real-time detection.

- It has existing open-source implementations that can be leveraged efficiently.

# Part 2: Implementation

To implement the **RawNet2** model for detecting AI-generated human speech, we'll proceed with the following steps:

## 1. Dataset Selection

We'll utilize the **ASVspoof 2019** dataset, a widely recognized benchmark in the field of audio deepfake detection. This dataset offers a comprehensive collection of both genuine and spoofed speech samples, making it ideal for training and evaluating our model.

## 2. Code Acquisition and Environment Setup

We'll leverage the open-source implementation of RawNet2 provided by the ASVspoof community. The repository is available on GitHub: rawnet2-antispoofing.

## 3. Data Preprocessing

Prepare the ASVspoof 2019 dataset:

- **Organize Data:** Structure the dataset into appropriate directories, typically train, dev, and eval, as expected by the RawNet2 implementation.

- **Feature Extraction:** RawNet2 processes raw audio waveforms directly, eliminating the need for handcrafted feature extraction. Ensure that audio files are correctly formatted (e.g., consistent sampling rates) as required by the model.

## 4. Model Training and Fine-Tuning

- **Initiate Training:** Execute the training script
- **Validation:** After training, evaluate the model on the development set to assess its performance. Utilize metrics such as Equal Error Rate (EER) and Detection Cost Function (DCF) for evaluation.

## 5. Implementation Comparison

**RawNet2 vs. AASIST:**

- *Architecture:* RawNet2 processes raw audio waveforms using Sinc-based convolutional layers, whereas AASIST employs spectro-temporal graph attention networks to capture intricate patterns in audio data.

- *Feature Extraction:* RawNet2's end-to-end approach eliminates the need for manual feature extraction, while AASIST relies on graph-based representations of spectrograms.

- *Complexity:* RawNet2's architecture is relatively lightweight, facilitating real-time applications. AASIST's graph-based approach may introduce additional computational overhead.

**RawNet2 vs. Dual-Branch Network:**

- *Architecture:* The Dual-Branch Network integrates both LFCC and CQT features, processing them through separate branches before fusion. RawNet2 utilizes a single stream, processing raw waveforms directly.

- *Feature Integration:* RawNet2 learns feature representations directly from data, whereas the Dual-Branch Network explicitly combines handcrafted features.

- *Performance:* Both models have demonstrated strong performance in detecting synthetic speech, with differences contingent on specific datasets and attack types.

# Part 3: Documentation & Analysis

# 1. Implementation Process

## Challenges Encountered

- **Dataset Format Handling**: The dataset was stored in Parquet format, which required proper loading and transformation into PyTorch tensors.
- **Computational Resources**: Training a deep learning model on large audio files required optimizing data loading and leveraging Google Colab's GPU.
- **Preprocessing Raw Audio**: Unlike standard spectrogram-based methods, RawNet2 operates on raw waveforms, necessitating careful data normalization and batch processing.

## Solutions Implemented

- Applied `torchaudio.transforms.Resample` to standardize sample rates.
- Implemented `DataLoader` with `num_workers` to speed up data loading.

## Assumptions Made

- The dataset labels were correctly assigned for training and evaluation.
- All audio samples had a consistent sampling rate after resampling.
- The model would generalize well to unseen real-world audio deepfake samples.

# 2. Analysis

## Model Selection: Why RawNet2?

1. **End-to-End Processing**: Unlike traditional methods that require handcrafted feature extraction, RawNet2 directly processes raw waveforms.
2. **Efficiency with Sinc Filters**: Uses Sinc-based convolution instead of standard CNNs, reducing computational overhead.
3. **Proven Performance**: RawNet2 has demonstrated strong results in past research on ASVspoof challenges.

## High-Level Technical Explanation

- **SincConv Layer**: Extracts frequency-related information using parameterized sinc functions.
- **Residual Blocks**: Helps learn deeper hierarchical representations.
- **LSTM Layer**: Captures temporal dependencies in speech patterns.
- **Fully Connected Layer**: Outputs probabilities for spoof vs. genuine classification.

## Performance Results

- **Dataset Used**: ASVspoof 2019 LA
- **Training Accuracy**: 99.71% (after 100 epochs)

## Strengths and Weaknesses

### Strengths

✔ End-to-end approach simplifies preprocessing.

✔ Efficient inference due to SincConv.

✔ Works well for raw waveform data.

### Weaknesses

❌ Requires significant training data to generalize well.

❌ Performance may drop in noisy environments.

❌ Not extensively tested for real-time deployment.

## Future Improvements

- Fine-tune on diverse, real-world datasets.
- Implement noise augmentation for robustness.
- Optimize for low-latency inference using quantization techniques.

# 3. Reflection Questions

## 1. Most Significant Challenges?

- Handling large audio files efficiently without memory bottlenecks.
- Fine-tuning hyperparameters for better generalization.
- Managing computational constraints in a cloud-based environment.

## 2. Real-World vs. Research Dataset Performance?

- Research datasets are curated and may not fully capture real-world variations (e.g., background noise, different recording devices).
- Real-world performance might degrade due to unseen attack strategies in deepfake generation.

## 3. Additional Data or Resources to Improve Performance?

- Collecting adversarial samples generated using the latest voice synthesis models (e.g., VALL-E, ElevenLabs).
- Augmenting data with noise, reverb, and different speech patterns.

- Leveraging semi-supervised learning to enhance performance with limited labeled data.

## 4. Deployment in Production?

- Use ONNX/TensorRT to optimize inference speed.
- Deploy as an API with a real-time streaming pipeline.
- Continuously monitor model drift and retrain with new data.
- Implement explainability techniques to improve trust in detection results.

# Conclusion

RawNet2, with its ability to process raw waveforms using SincConv and LSTMs, offers a promising approach for audio deepfake detection. While it performed well on ASVspoof 2019, further improvements in real-world robustness and deployment efficiency are necessary. Future work should focus on improving generalization through adversarial training and optimizing inference for real-time applications.

**Code**: [Colab notebook](#)
**Dataset**: [Dataset](#)
**Github**: [Repo](#)