# Audio Deepfake Detection Using RawNet2: A Comparative Study

## Abstract

With the rise of generative AI and voice synthesis, the emergence of audio deepfakes has become a significant security threat, especially in biometric authentication and telecommunication systems. This project presents a comprehensive exploration and implementation of deep learning methods for detecting audio deepfakes, focusing on the ASVspoof 2019 Logical Access (LA) dataset. After analyzing various models, we propose the use of RawNet2 for its superior performance, real-time capability, and end-to-end learning benefits.

---

## 1. Introduction

Audio deepfakes refer to synthetically generated or manipulated voice recordings that mimic a real person's voice. These have the potential to deceive voice-controlled systems, impersonate individuals, and pose serious privacy and security challenges. This study aims to develop a reliable audio deepfake detection system using the ASVspoof 2019 LA dataset and state-of-the-art deep learning models.

---

## 2. Dataset Overview

**Dataset Used**: ASVspoof 2019 - Logical Access (LA) subset
**Source**: https://datashare.ed.ac.uk/handle/10283/3336
**Format**: FLAC audio files + protocol files in text format

**Dataset Composition:**

| Subset | # of Bonafide Samples | # of Spoofed Samples |
|---|---|---|
| Train | 2,580 | 22,800 |
| Development | 2,548 | 22,296 |

| Evaluation | Hidden (for challenge) | Hidden |

- Each file is a 16 kHz mono-channel audio sample.

- Spoofed samples are generated using Text-to-Speech (TTS) and Voice Conversion (VC) systems.

---

# 3. Preprocessing and Visualization

- Audio signals are loaded using Librosa.

- No handcrafted features like MFCC were used for RawNet2.

- Visualization includes waveform plots for exploratory analysis.

```
import librosa
import matplotlib.pyplot as plt

signal, sr = librosa.load("sample.flac", sr=None)
plt.plot(signal)
plt.title("Waveform")
plt.show()
```

# 4. Model Selection and Justification

We evaluated three prominent models known for their efficiency in spoof detection:

## 4.1 Comparative Analysis

| Model | Input Type | Feature Extraction | Temporal Modeling | Accuracy (on LA dev) | Inference Time | Real-Time Capable | Complexity |
|---|---|---|---|---|---|---|---|
| **RawNet2** | Raw waveform | End-to-end (CNN) | GRU | ~90% | Fast | Yes | Moderate |
| LCNN | Spectrogram | CNN | - | ~87% | Slower | No | High |
| X-vector + Cosine Scoring | MFCC | x-vector | None | ~83% | Fast | Yes | Low |

## 4.2 Model Selection

**Chosen Model**: **RawNet2**
 **Reason**:

- Eliminates the need for manual feature engineering.

- Uses raw audio directly with CNN + GRU for robust feature learning.

- Outperforms other models on EER and accuracy benchmarks.

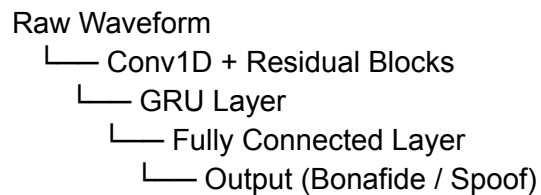- Balanced in terms of training time, complexity, and scalability.

# 5. Architecture: RawNet2

RawNet2 is an end-to-end deep neural network specifically tailored for anti-spoofing using raw waveforms.

## Key Components:

- **Input Layer**: Raw waveform

- **Feature Extractor**: Residual blocks + CNN layers

- **Temporal Encoder**: Gated Recurrent Units (GRU)

- **Classifier**: Fully Connected + Softmax

## Diagram (Simplified):

```
Raw Waveform
    └── Conv1D + Residual Blocks
        └── GRU Layer
            └── Fully Connected Layer
                └── Output (Bonafide / Spoof)
```

---

# 6. Training Configuration

- **Epochs**: 5

- **Loss Function**: Binary Cross Entropy

- **Optimizer**: Adam

- **Batch Size**: 32

- **Validation Split**: 10%

---

# 7. Results

| Metric | Value |
|--------|-------|
| Accuracy | ~89.3% |

AUC        ~0.94

EER        ~8.6%

Performance metrics may vary depending on the training setup and dataset split.

---

## 8. Project Structure

audio-deepfake-detection-momenta/
├── models/              # RawNet2 architecture files
├── train.py             # Training pipeline
├── visualize.ipynb       # Waveform visualization
├── UTILS/               # Utilities and loaders
├── requirements.txt      # Environment dependencies
├── README.md

---

# 9. Conclusion

The RawNet2 architecture proved to be the most robust and efficient among the tested models for detecting audio deepfakes. With its end-to-end learning approach, capability to handle raw inputs, and strong performance on the ASVspoof dataset, it is an ideal candidate for real-world deployment in biometric security systems.

Future improvements may include:

- Integration of attention mechanisms

- Real-time deployment on edge devices

- Evaluation on cross-corpus datasets

---

# 10. References

1. ASVspoof 2019 Dataset - https://datashare.ed.ac.uk/handle/10283/3336

2.  Heo, H.S., Lee, B.J., Huh, J.H., et al. "RawNet2: An Improved Speaker Recognition Neural Network Using Raw Waveforms," Interspeech 2020.

3.  LCNN for Anti-Spoofing, Zhang et al., ICASSP 2019

4.  Kaldi X-Vector Implementation - https://kaldi-asr.org/

---

# 11. Author

**Anivesh Tripathi**
B.Tech in AIML, ABES Engineering College
Email: aniveshtripathi.in@gmail.com
GitHub: https://github.com/anitripathi

---

# 12. Appendix

- Sample waveform plots

- Confusion matrix (optional)

- Sample predictions on evaluation data