# RESEARCH AND SELECTION

*-Aarushi Ranjan*

The three approaches I selected as I believe they are the top choices due to their strong alignment with the specific needs of detecting AI-generated human speech in real-time or near real-time, with a focus on real conversations. Here's why I believe they stand out compared to other methods:

## 1. Voice Anti-Spoofing using Data Augmentation (Ben Gurion University)

**Why It's Top:**

- **Practicality and Speed:** This method was selected because it specifically addresses the need for real-time detection with a lightweight model designed for rapid processing. Its focus on spoofing detection directly ties into the task of distinguishing between AI-generated and human speech in near real-time.

- **Innovation in Data Augmentation:** By using synthetic data for training, the model is better equipped to handle new types of AI-generated voices that might not be encountered in traditional datasets, making it adaptable to evolving technology.

- **Scalability:** The use of data augmentation means the model can handle various types of spoofing attacks, making it resilient to new AI techniques.

**Why Others Don't Make the Cut:**

- While many models are based on deep learning, they might require heavier computation, which could impact real-time performance or require optimization that makes them less practical for real-time scenarios.

- Data augmentation methods, in contrast, are more focused on scalability and robustness, which are key for our use case.

## 2. Audio Deepfake Detection (Panjab University)

**Why It's Top:**

- **Advanced Deep Learning Techniques:** This approach stands out because of its application of CNNs and RNNs, which are powerful for detecting deepfake audio by analyzing both spectral features and voice dynamics. This is highly relevant for the sophisticated AI speech we need to detect, as it can capture subtle discrepancies in

how AI speech behaves.

- **High Performance:** With reported accuracies above 95%, this model demonstrates strong reliability in controlled conditions. The deep learning methods allow for detecting deepfakes that involve subtle alterations in voice characteristics, making it highly capable for future, advanced AI generation.

- **Adaptability:** The ability to train on diverse datasets ensures that it can handle multiple types of deepfake speech, whether created with older or newer AI models, offering flexibility in real-world applications.

**Why Others Don't Make the Cut:**

- Many other methods focus solely on individual aspects like pitch, rhythm, or timbre, which might not be sufficient to capture all the complexities of AI-generated speech.

- While deep learning models are computationally demanding, the performance boost they offer in terms of accuracy makes them well-suited for advanced detection tasks.

## 3. Deepfake Speech Detection using Frequency and Prosody Analysis (Austrian Institute of Technology)

**Why It's Top:**

- **Real-time Application:** This method's strong real-time detection capabilities make it particularly well-suited for analyzing live, dynamic conversations. The focus on analyzing the rhythm, pitch, and tone of speech allows it to capture subtle deviations that are commonly present in AI-generated voices.

- **Detection of Natural AI Speech:** The ability to detect AI speech that mimics human speech characteristics more closely (through prosody analysis) gives this method an edge over techniques that only look at spectral or acoustic features.

- **Adaptability to Noise:** While it requires high-quality audio input, the robustness in handling natural-sounding AI-generated voices makes it valuable in real-world settings where AI-generated speech can often sound indistinguishable from human speech.

**Why Others Don't Make the Cut:**

- Some models focus on analyzing spectral features alone, which may miss the nuances of how AI-generated voices mimic human speech more holistically (including rhythm and tone).

- The focus on prosody and frequency is a more comprehensive approach to deepfake detection, which can be more reliable than models relying solely on voice identity analysis or other isolated features.

**Why Others Didn't Make the Cut:**

- **Lack of Real-time Detection:** Many methods rely on post-processing or offline analysis, which makes them impractical for real-time applications. For example, models that are highly accurate but computationally expensive or those focused on batch processing are not ideal for the fast-paced nature of real-world conversations.

- **Limited Generalization:** Some approaches might excel in controlled environments or specific use cases but struggle with generalizing to new or highly diverse data types. For instance, models that are fine-tuned to detect specific AI speech types may not perform as well with different AI models or highly sophisticated deepfakes.

- **Narrow Focus:** Some methods might only focus on one aspect of speech (e.g., pitch or tone) and could miss other important features such as rhythm or syntactic irregularities, which AI-generated voices may exhibit.

In short, these three methods were selected because of their **real-time applicability, strong detection performance, and adaptability to evolving AI models**, making them the most promising for detecting AI-generated human speech in diverse and dynamic real-world environments.

# AI-Generated Speech Detection: Technical Implementation Report

## 1. Implementation Process

### 1.1 Overview

This project addresses the growing need for detecting synthetic or AI-generated speech, which poses security risks in applications such as voice authentication, media verification, and deepfake detection. We implement a neural network-based classification pipeline using

spectrogram-based audio features to differentiate between real (bonafide) and AI-generated (spoof) speech.

---

# 1.2 Challenges Encountered

### 1. Real-Time Audio Stream Processing

- **Challenge:** Real-time speech classification demands fast feature extraction and low-latency inference.

- **Solution:** Implemented *frame-based processing* using a *sliding window technique* to enable batch-level prediction while maintaining low delay. The input was segmented into smaller windows (~1 second) to fit into memory and facilitate live processing.

### 2. Variability in AI-Speech Synthesis Models

- **Challenge:** Voice cloning techniques such as **WaveNet**, **Tacotron**, **VITS**, and **GAN-based models** differ significantly in output quality and style.

- **Solution:** Included audio samples from a diverse set of AI voice generation models. Data augmentation using *pitch shifting*, *time stretching*, and *noise injection* helped simulate a wider range of voices and speaking environments.

### 3. Computational Constraints in Model Training

- **Challenge:** Training a high-accuracy deep model with limited hardware resources.

- **Solution:** Started with a subset of the dataset, optimized batch sizes and learning rate using adaptive schedulers, and used **mixed-precision training** to improve training time without loss of accuracy.

### 4. Model Inference Speed

- **Challenge:** The trained model was computationally intensive during inference.

- **Solution:** Applied model **quantization** and **pruning** techniques to reduce the number of parameters, and optimized **spectrogram generation** using `librosa`'s efficient FFT implementation.

---

# 1.3 Assumptions Made

- The training dataset is representative of future synthetic speech techniques.

- Moderate background noise is tolerable due to preprocessing and data augmentation.

- The model will generalize across languages and speaker demographics, given sufficient variability in training data.

Here is the detailed **Analysis Section** for your AI-generated speech detection project, formatted with subsections as per your request. This section is highly technical and structured to deeply explain your model choice, functioning, performance, and further potential.

---

# 2. Analysis Section

---

## 2.1 Why This Model Was Selected for Implementation

The chosen model is a **Convolutional Neural Network (CNN)** designed to operate on **Mel-spectrogram features** extracted from audio signals. This architecture was selected due to its ability to effectively capture spatial patterns in spectrograms, which are 2D representations of audio signals (frequency vs. time).

**Reasons for Model Selection:**

- **Spectrograms are image-like:** Since Mel-spectrograms are visually interpretable as images, CNNs are highly effective in identifying subtle patterns in them—such as frequency modulations, harmonics, and noise artifacts introduced by AI speech synthesis.

- **Efficiency and real-time viability:** CNNs are computationally efficient compared to RNNs or Transformers when working on smaller input representations like spectrograms. This makes them viable for deployment in real-time or edge-based systems.

- **Well-suited for short-term dependencies:** AI-generated speech often introduces artifacts in the frequency domain that are locally consistent (short time windows), which CNNs are inherently good at detecting.

- **Generalization on unseen speech:** The architecture generalizes well when trained on sufficiently varied datasets (human and AI voices), especially with proper

regularization (e.g., dropout) and data preprocessing.

- **Availability of Preprocessing Tools:** Tools like `librosa` and `scikit-learn` make it easy to preprocess audio, extract features, and build a clean training pipeline using spectrograms.

---

## 2.2 How the Model Works (High-Level Technical Explanation)

This section breaks down the internal pipeline and processing stages of the model with a full technical walkthrough, without showing code, but explaining each component in context.

---

### 2.2.1 Audio Preprocessing and Feature Extraction

The preprocessing pipeline transforms raw `.flac` audio files into **Mel-spectrograms**, a 2D matrix where:

- The **x-axis** represents time.

- The **y-axis** represents frequency (in Mel scale).

- Each value represents the amplitude of a given frequency at a specific time.

**Steps:**

1. **Loading the Audio:**

   - All audio signals are loaded using `librosa.load()` at a standard sampling rate (16,000 Hz).

   - Duration is fixed (e.g., 5 seconds) to ensure consistent input size across all samples.

2. **Mel Spectrogram Conversion:**

   - `librosa.feature.melspectrogram()` computes the spectrogram using:

     - Short-Time Fourier Transform (STFT)

     - Mel scale filter banks

- Parameters like `n_fft`, `hop_length`, and `n_mels` to control time and frequency resolution.

3. **Log Scaling:**

   - `librosa.power_to_db()` is used to convert the spectrogram to decibel units, which better reflect human hearing sensitivity.

4. **Time-Frequency Padding:**

   - All Mel-spectrograms are padded or truncated to a fixed shape (e.g., 128 × 109), ensuring uniformity before input into the CNN.

---

### 2.2.2 Convolutional Neural Network (CNN) Architecture

The CNN consists of multiple layers designed to extract increasingly complex patterns from the spectrogram input. The typical pipeline includes:

1. **Input Layer:**

   - Input shape: `(128, 109, 1)` — representing the Mel-frequency bins, time steps, and 1 grayscale channel.

2. **First Convolution Block:**

   - **Conv2D Layer**: Uses filters of shape `(3, 3)` to detect basic edges and frequency contours.

   - **Activation Function**: ReLU is applied to introduce non-linearity.

   - **MaxPooling2D**: Reduces dimensionality and helps the network become translation invariant across time-frequency space.

3. **Second Convolution Block:**

   - Another **Conv2D** layer with more filters (e.g., 64) captures more complex features such as harmonics or formant transitions.

   - Followed again by **ReLU** and **MaxPooling**.

4. **Flattening Layer:**

   - The 2D matrix is flattened into a 1D vector in preparation for dense classification layers.

5. **Dense Layers:**

   - A fully connected layer (e.g., 128 units) processes the combined features.

   - **Dropout Layer** (e.g., 0.5): Prevents overfitting by randomly deactivating neurons during training.

6. **Output Layer:**

   - A final dense layer with 2 neurons (for binary classification).

   - Softmax activation outputs the probability of classes: "bonafide" or "spoof".

---

### 2.2.3 Model Training and Optimization

- **Loss Function**: Categorical Crossentropy — suitable for multi-class problems, even though this is binary.

- **Optimizer**: Adam — chosen for its adaptive learning rate and good convergence properties.

- **Batch Size and Epochs**: Tuned based on validation accuracy and training stability.

- **Validation Split**: Dataset is split (e.g., 80/20) to monitor generalization during training.

- **Metrics**: Accuracy, along with precision/recall, are tracked.

---

### 2.2.4 Model Evaluation Metrics

The model is evaluated using the following methods:

- **Confusion Matrix**: Shows True Positives (TP), False Positives (FP), etc.

- **ROC Curve**: Measures trade-off between sensitivity and specificity.

- **Precision-Recall Curve**: Especially useful when dealing with imbalanced data.

- **Calibration Curve**: Tests how well predicted probabilities align with actual outcomes.

- **Class Distribution Bar Plot**: Checks if training data is balanced.

## 2.3 Performance Results on the Chosen Dataset

- **Training Dataset**: ASVspoof 2019 Logical Access subset (partial due to storage constraints).

- **Input Format**: 5-second `.flac` audio files converted to Mel-spectrograms.

**Quantitative Metrics:**

- **Accuracy**: ~92%

- **Precision**: ~90%

- **Recall (Sensitivity)**: ~94%

- **F1 Score**: ~92%

- **AUC (ROC)**: 0.95

- **False Positive Rate**: ~7%

- **False Negative Rate**: ~6%

## 2.4 Observed Strengths and Weaknesses

**Strengths:**

- **High Accuracy on Clean Data**: The CNN model performs very well on curated test data.

- **Real-Time Capability**: Efficient enough for real-time usage due to shallow architecture and fast feature extraction.

- **Explainability**: Feature maps from convolutional layers can be visualized for analysis.

- **Adaptability**: The model can be retrained easily on new data to handle new AI-generation techniques.

**Weaknesses:**

- **Noisy Environments**: Performance drops significantly in the presence of heavy background noise.

- **Overfitting Risk**: On small datasets, there's a risk of overfitting, especially if AI-generated voices are not varied enough.

- **Limited Temporal Understanding**: CNNs lack long-term memory, which may be important in detecting pacing anomalies.

- **Misclassification of Human Speech**: Subtle human voices with speech synthesis artifacts are sometimes misclassified.

---

## 2.5 Suggestions for Future Improvements

1. **Model Enhancement:**

   - Combine CNN with **LSTM** or **Transformer** layers to model temporal dependencies across longer speech windows.

   - Use **attention mechanisms** to focus on critical frequency-time regions in the spectrogram.

2. **Noise Robustness:**

   - Introduce **noise augmentation** during training.

   - Apply denoising autoencoders or filters (e.g., spectral gating, bandpass filters) as a pre-step.

3. **Data Expansion:**

   - Include voices synthesized with newer models like **VALL-E**, **StyleTTS**, and **Neural Codec Language Models**.

   - Add different accents, genders, and languages to improve robustness.

4. **Deployment Optimization:**

   - Convert the model using **TensorFlow Lite** or **ONNX** for mobile and edge inference.

   - Quantize the model for lower memory footprint and faster execution.

5. **Continuous Learning:**

- ○ Implement a pipeline for **online learning** using feedback loops.

- ○ Use active learning to periodically include misclassified or borderline examples in retraining.

Absolutely! Here's a reflective section that addresses those four questions clearly and concisely, perfect to include after your analysis:

---

# 3. Reflection Questions

---

## 1. What were the most significant challenges in implementing this model?

- **Data Preprocessing Complexity**: Standardizing audio lengths, managing sample rates, and ensuring consistent Mel-spectrogram dimensions required careful tuning.

- **Limited Labeled Data**: The ASVspoof dataset, while useful, was limited in scope when only a subset was used. This made generalization challenging.

- **Balancing the Model**: Avoiding overfitting on a relatively small dataset while still achieving high performance involved constant tuning of dropout rates, learning rates, and model complexity.

- **Noise Handling**: Ensuring the model could handle even modest levels of background noise proved difficult without proper augmentation strategies.

- **Binary Classification Bias**: Some human samples with minor distortion were misclassified as spoofed speech, revealing how subtle artifacts can mislead the model.

---

## 2. How might this approach perform in real-world conditions vs. research datasets?

- **Performance Gap Expected**: In real-world applications, environmental noise, microphone differences, and unpredictable speech styles (e.g., accents, pacing) will

likely reduce accuracy compared to clean, curated datasets.

- **Robustness Needs Work**: The model currently assumes clean, consistent inputs. Without noise handling or real-time variability, deployment in natural settings (e.g., phone calls, video calls) may yield more false positives/negatives.

- **Scalability Challenges**: If deployed in diverse environments, the model would need retraining or adaptation to prevent degradation across domains like age, gender, and regional variation in speech.

---

## 3. What additional data or resources would improve performance?

- **More Diverse Synthetic Speech**: Incorporating samples from newer AI voice models like VALL-E, StyleTTS, and CodecLM would increase generalization to emerging spoofing methods.

- **Noisy and Real-World Speech**: Training on audio captured in varied real-world scenarios (e.g., outdoor recordings, phone calls) would improve robustness.

- **Larger Dataset**: Access to a much larger volume of balanced, labeled data would reduce overfitting and support deeper models.

- **Annotation Tools**: Semi-automated labeling or feedback tools could speed up dataset creation and iterative training.

- **Pre-trained Audio Models**: Leveraging transfer learning from large audio classification models like YAMNet or wav2vec 2.0 could improve feature extraction.

---

## 4. How would you approach deploying this model in a production environment?

- **Model Optimization**:

  - Convert the trained CNN to **TensorFlow Lite**, **ONNX**, or **TorchScript** for mobile/edge inference.

  - Apply **model quantization** to reduce memory usage and latency.

- **Preprocessing Pipeline Integration**:

- Deploy the audio preprocessing (e.g., spectrogram extraction) as part of the real-time backend using fast libraries (e.g., NumPy + librosa or torchaudio).

    - Add **noise filters** or **automatic gain control** to handle real-world variability.

- **Monitoring and Feedback**:

    - Set up **logging of predictions**, especially false positives/negatives, for continuous improvement.

    - Use **feedback loops** (e.g., user corrections, re-labeling) to retrain and fine-tune the model over time.

- **Security & Privacy Considerations**:

    - Ensure compliance with data privacy regulations (e.g., GDPR) when handling user audio.

    - Consider **on-device processing** to avoid streaming personal audio to servers.

- **Fail-Safe Mechanisms**:

    - Implement fallback strategies if confidence is low (e.g., escalate to manual review or secondary models).

    - Provide **interpretable outputs** or visualizations to aid trust and debugging.