

WaveTruth

Cracking Deepfakes with Deep Learning



What is WaveTruth

WaveTruth is an advanced web application designed to analyze audio recordings and determine:

- ✅ AI vs. Real – Detects whether a voice is generated by artificial intelligence or a real human.
- ✅ Gender Identification – Classifies the speaker as male or female based on voice characteristics.



How It Works:

Powered by a Feedforward Neural Network trained through transfer learning, WaveTruth ensures high-accuracy detection by using deep learning techniques in audio analysis.



Bringing Transparency to the Age of AI Voice Generation!

Scope and Audience

Scope

WaveTruth is designed to classify audio clips as either AI-generated or real human voices by analyzing key audio features such as:

- 🎵 Pitch Variability – Detects unnatural pitch fluctuations
- 🔊 Spectral Patterns – Identifies AI-generated sound signatures
- 💡 Prosody Analysis – Examines rhythm, tone, and speech dynamics

To enhance accuracy, the model is trained on a diverse dataset of AI and human voices, including varied accents, tones, pitch, and speech speed.

Audience

WaveTruth serves a broad range of users, including:

- 👤 Content Moderators – Detects AI-generated voices in media
- 🛡️ Security & Fraud Prevention Experts – Identifies synthetic speech in scams and deepfakes
- 🔬 AI & Voice Recognition Researchers – Advances studies in audio authentication

🚀 Bringing Trust & Transparency to Voice AI!

Background Review

Existing Approaches

Several companies have developed advanced solutions for detecting AI-generated voices by utilizing cutting-edge machine learning and deep learning techniques. For example, Evenlabs uses a variety of AI models, including Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), to detect synthetic speech. This approach allows Evenlabs to accurately identify synthetic voices, but the system is limited by access restrictions, high costs, and challenges with lower-quality or heavily edited audio. PlayHT, on the other hand, combines neural networks, including Generative Adversarial Networks (GANs) and RNNs, for both voice synthesis and detection. While PlayHT provides high-quality detection and synthesis, it also faces limitations such as potential issues with noisy audio and non-transparent pricing structures. Resemblyzer leverages wav2vec 2.0 embeddings for speaker verification and similarity detection, excelling at distinguishing specific speakers based on their voice features. However, its limitations include struggles with regional accents and languages with fewer resources.

While these companies offer innovative solutions for detecting synthetic voices, they each face unique challenges in terms of generalization across diverse audio environments and the evolving sophistication of AI voice generation techniques. As the field advances, there is a growing need for models that integrate multiple audio patterns and improve detection accuracy, which is the primary goal of this project.

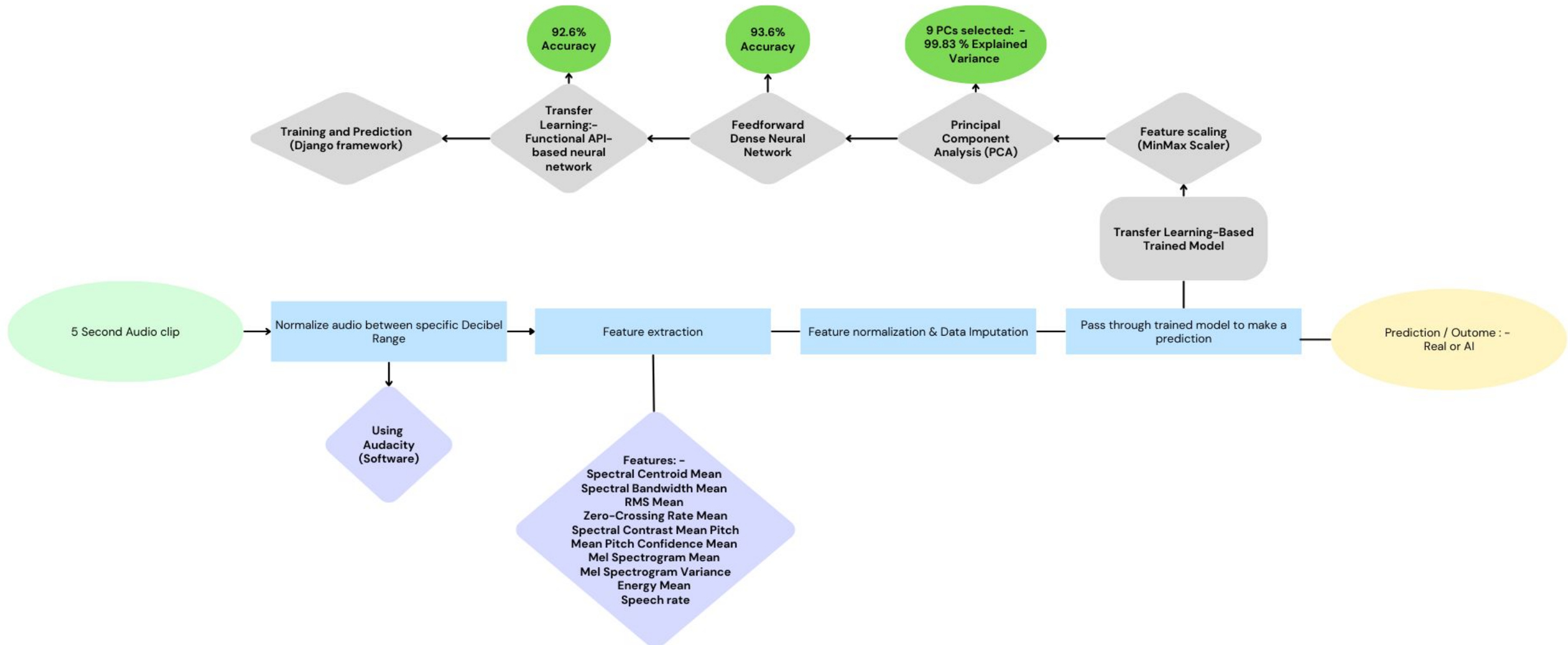
Related Literature

The exploration of voice synthesis and detection has evolved significantly over the years, with early work focusing on techniques like Mel Frequency Cepstral Coefficients (MFCCs). Recent research has increasingly turned to deep learning models such as CNNs and RNNs to classify and differentiate between human and synthetic voices. However, even these models often struggle to accurately identify new or previously unseen synthetic voices. Transformer-based architectures, which enhance context-awareness for processing audio data, have recently been introduced and show promise in improving detection accuracy. Building on this body of work, this project proposes a hybrid model that combines CNNs and transformers to better capture both temporal and spatial audio patterns, offering an improved approach to synthetic voice detection.

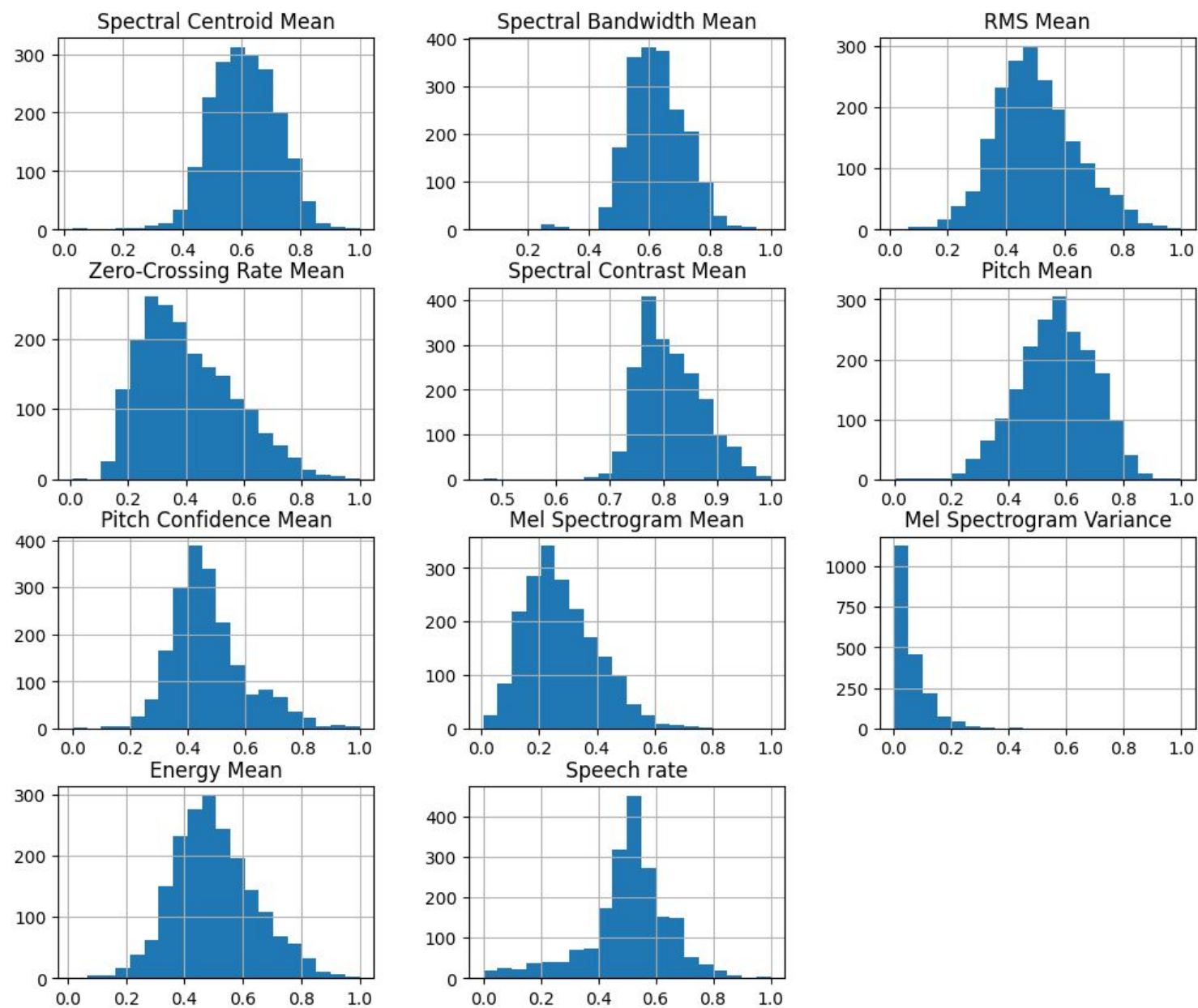
WaveTruth vs. Other Similar AI Detectors: A Comparative Analysis

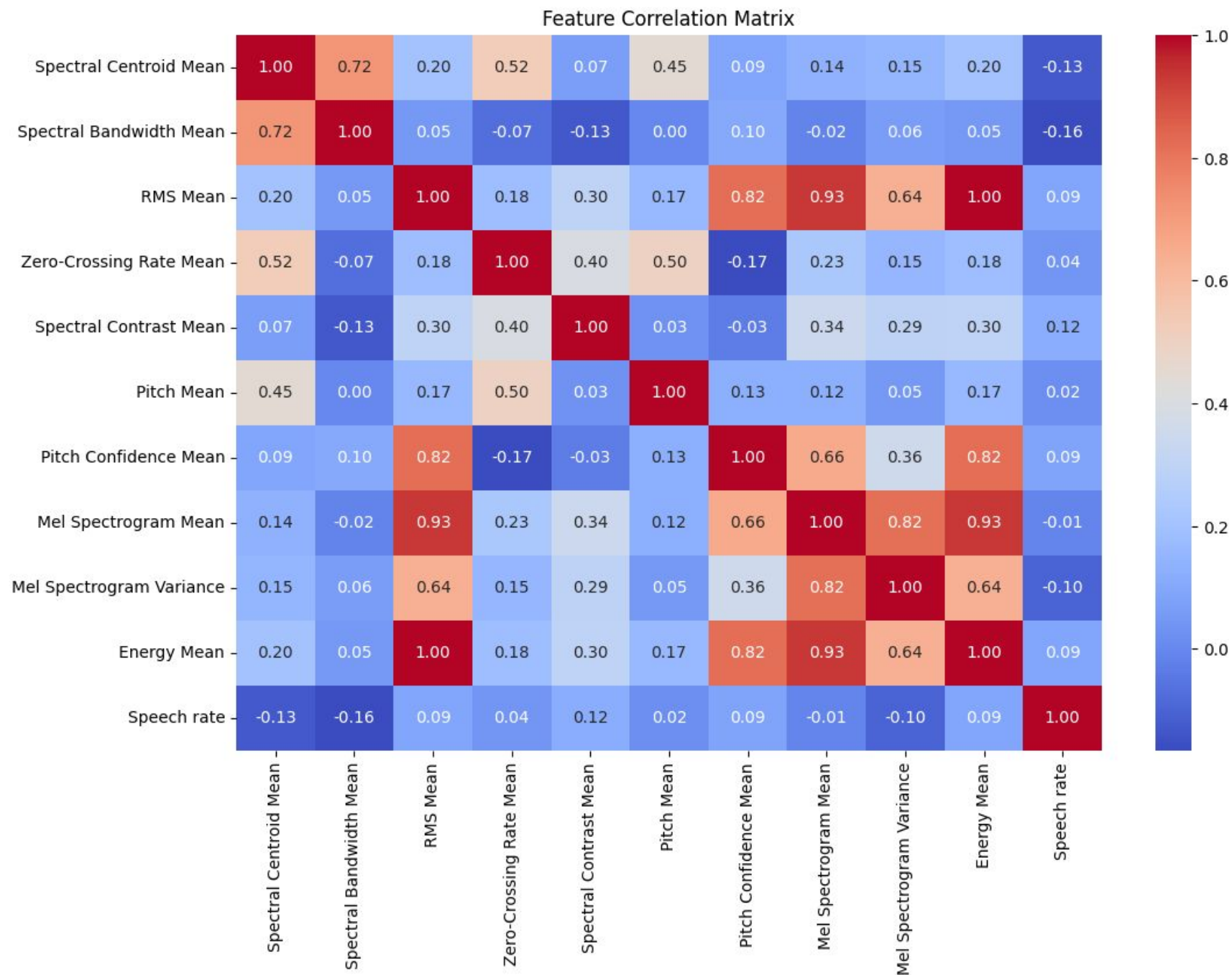
Feature	Evenlabs	PlayHT	Resemblyzer	WaveTruth
Detection Model/Algorithm	Uses a variety of AI models for synthetic speech detection, including deep learning techniques like CNNs and RNNs	Uses a combination of neural networks (including deep learning models like GANs and RNNs) for voice synthesis and detection	wav2vec 2.0 embeddings for speaker verification and similarity detection	<i>Feed-forward neural networks with transfer learning for gender classification</i>
Drawbacks	<ul style="list-style-type: none">- Limited access; enterprise-focused.- Can struggle with lower-quality or heavily edited audio.- Requires commercial licensing.	<ul style="list-style-type: none">- High cost for high-volume use.- May struggle with noisy or low-quality audio.- Pricing based on usage, not always transparent.	<ul style="list-style-type: none">- False positives for similar-sounding speakers.- Struggles with regional accents and low-resource languages.- Limited access, mainly for research or enterprise use.	<ul style="list-style-type: none">- <i>Limited by dataset size for underrepresented voices, though transfer learning helps mitigate this.</i>- <i>May struggle with synthetic-sounding voices.</i>- <i>Accessible via a web app only currently</i>
Strengths	<ul style="list-style-type: none">- High accuracy in synthetic speech detection.- Well-suited for enterprise-level use and security applications.	<ul style="list-style-type: none">- High-quality synthetic speech generation and detection.- Great for personalized voice creation.	<ul style="list-style-type: none">- Efficient at speaker verification and comparison.- Effective for distinguishing specific speakers based on audio features.	<ul style="list-style-type: none">- <i>Focuses on gender classification with a specialized and accessible approach.</i>- <i>Uses transfer learning to improve generalization on limited datasets.</i>- <i>Free to use</i>
Goal/Focus	To provide cutting-edge solutions for detecting synthetic voices and deepfakes in audio.	To offer high-quality text-to-speech synthesis and detection tools, including for AI-generated voices.	To verify and differentiate speakers based on their voice characteristics.	<i>To provide solutions for detecting AI audio and make gender-based voice classification accessible and reliable.</i>

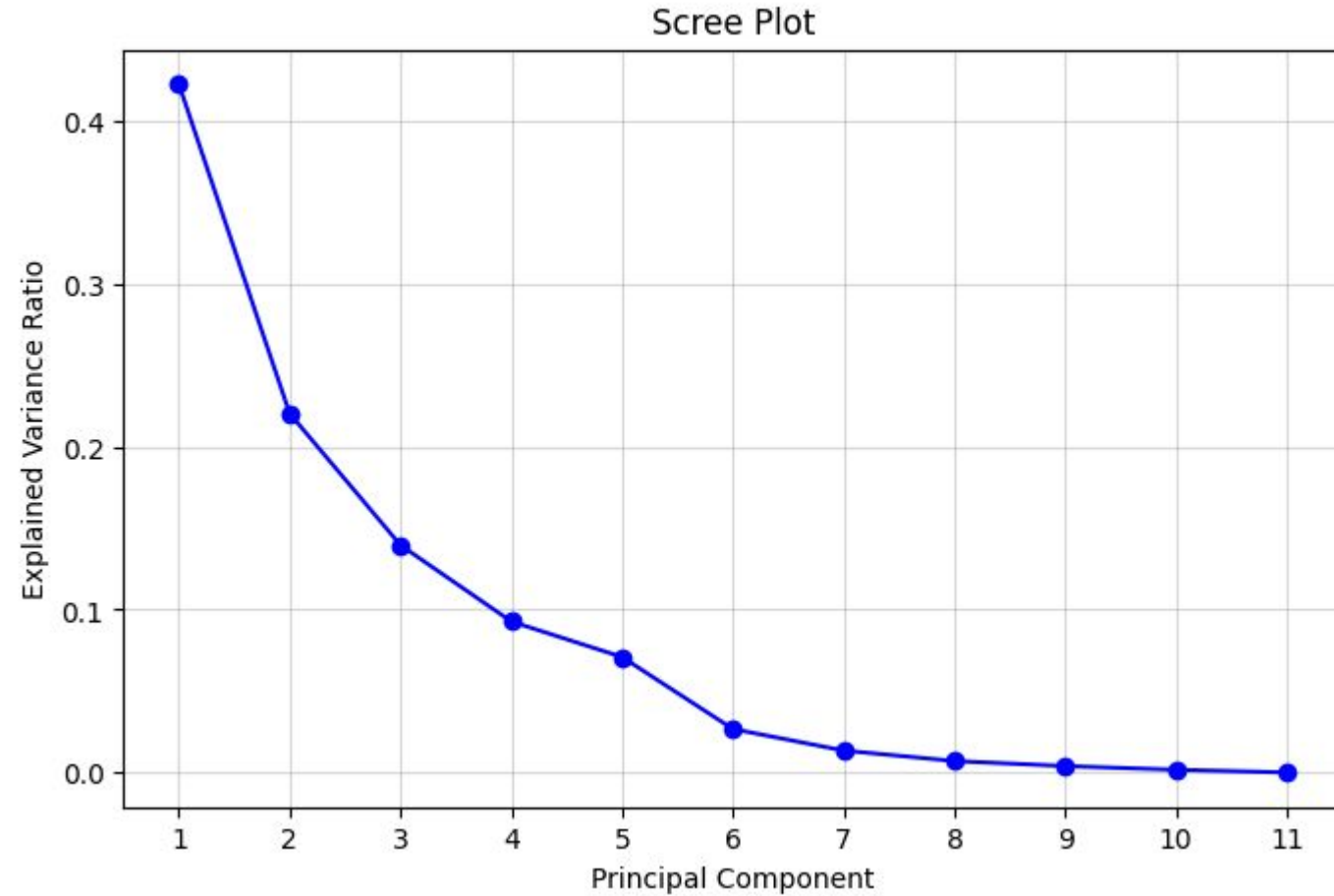
Pipeline: - Behind the scenes of how the WaveTruth makes its predictions



Feature Distributions

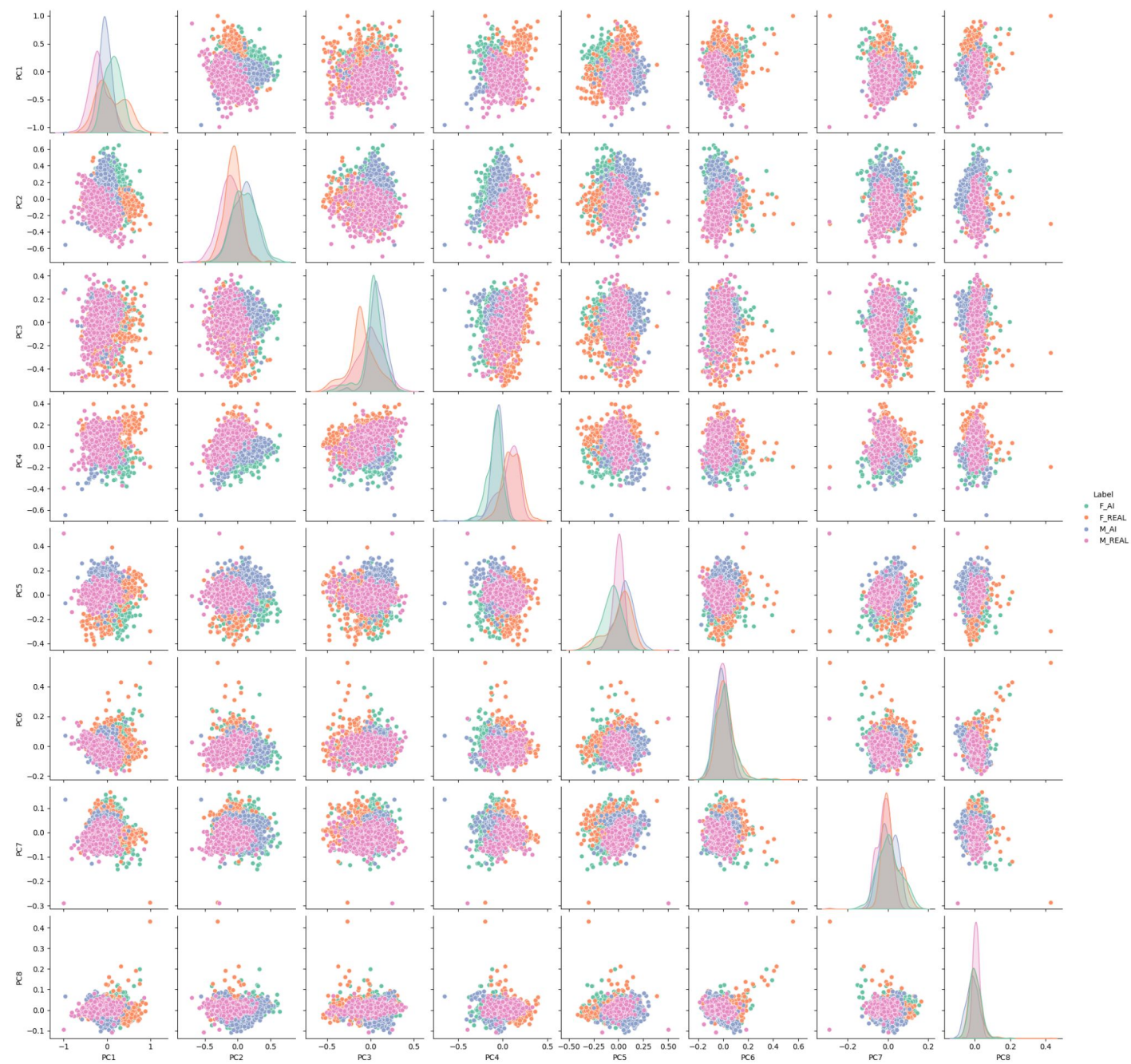


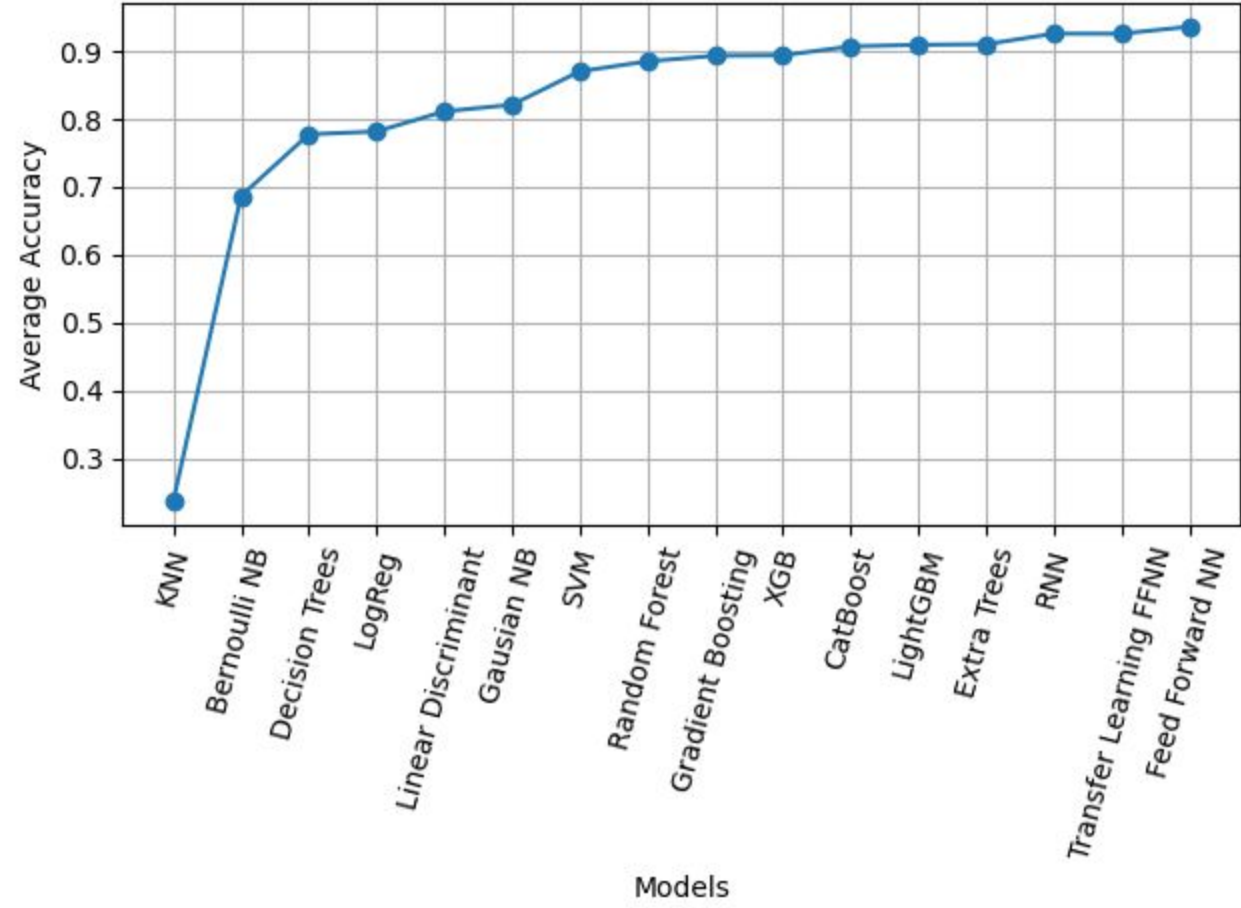




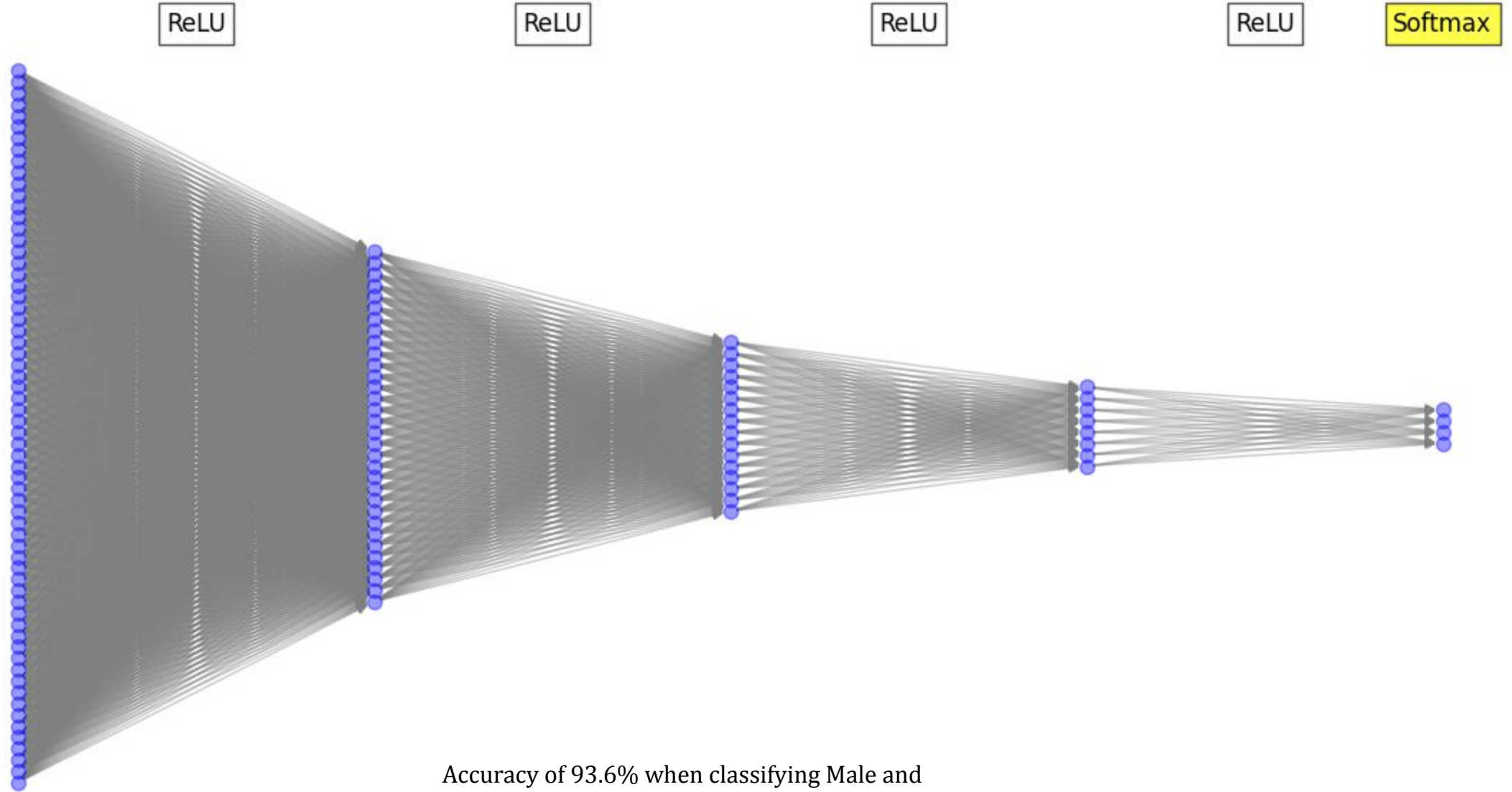
9 Principal Components Selected, 99.83 %
of total explained Variance

Pairwise Scatter Plots of Principal Components



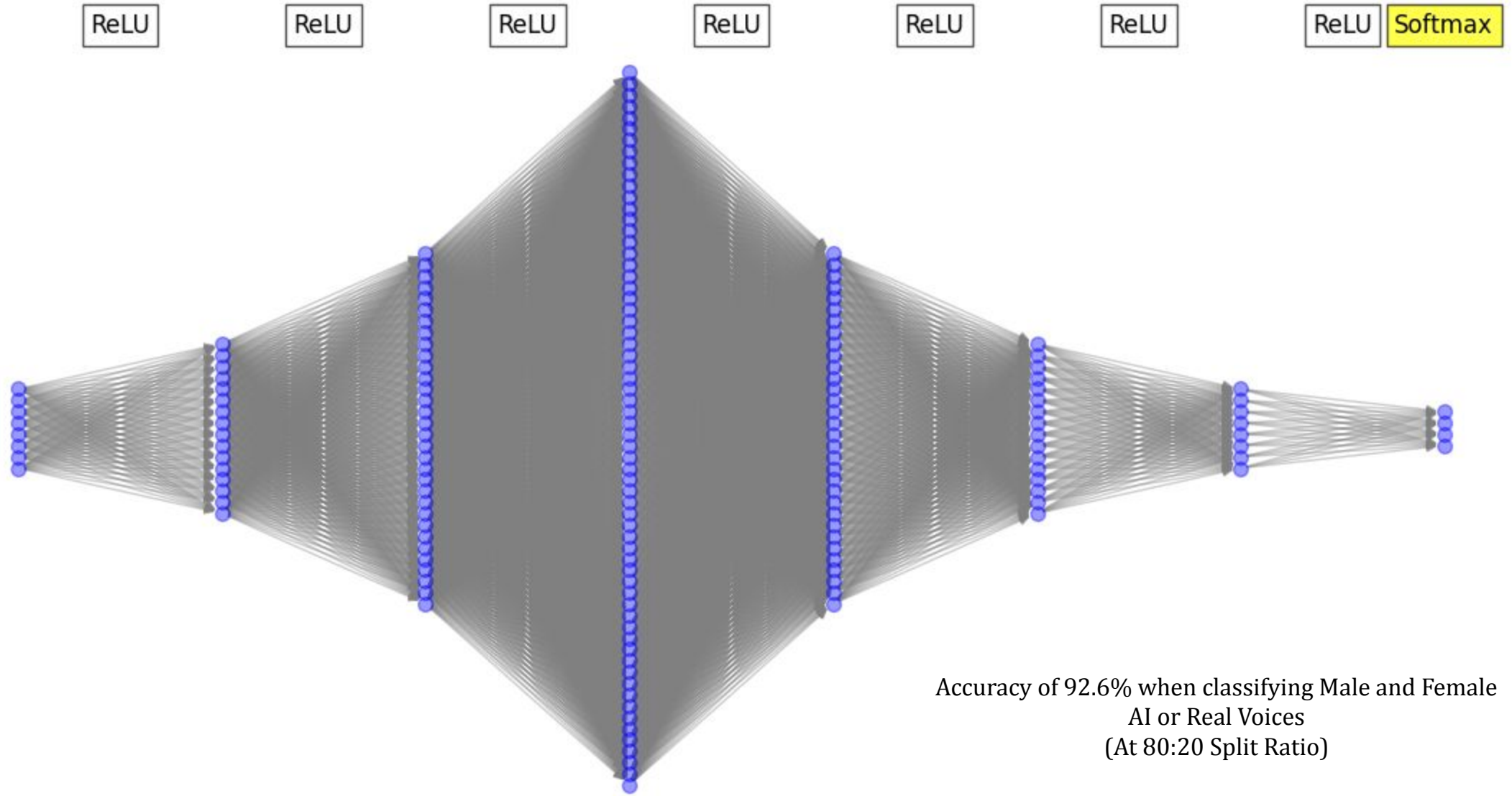


Feed Forward Neural Network Visualization



Accuracy of 93.6% when classifying Male and Female AI or Real Voices
(At 80:20 split ratio)

Transfer Learning Neural Network Visualization



UI development in progress

Django Framework