

FakeVoiceGuard: A Hybrid ResNeXt BiGRU Transformer Framework for Robust Deepfake Audio Detection Using ASVspoof Datasets

R Ranjith¹, C Viswanathan², B Tamil Selvan³

^{1,2,3}*Department of Computer Science and Engineering, SRG Engineering College, Namakkal, Tamil Nadu*

Abstract—The rise of deepfake audio technologies that can create highly realistic synthetic voices has, among other things, digital forensic, biometric authentication, and media integrity challenges. The detection of such faked voices calls for models that can perfectly grasp the spatial, temporal, and contextual nature of the speech signals. The present work introduces FakeVoiceGuard, a hybrid deep learning model that combines ResNeXt, Bidirectional Gated Recurrent Unit (Bi-GRU), and Transformer-based classification to pinpoint artificial or spoofed voices with a high degree of accuracy. In the newly designed method, to capture the time and frequency characteristics, input audio samples in the form of log-mel spectrograms are generated. The ResNeXt encoder, through grouped residual convolutions, extracts the deep spectral features, while the Bi-GRU layer captures the bidirectional temporal dependencies of the speech. At last, the Transformer unit, employing multi-head self-attention, determines the most informative temporal segments for the classification. We test the model on the ASVspoof 2019 (Logical and Physical Access) and ASVspoof 2021 (LA, PA, and Deepfake) datasets, which contain a variety of spoofing attacks such as text to speech synthesis, voice conversion, and replay audio, and subsequently train it on them. Experimental results point to the fact that FakeVoiceGuard is the most accurate and robust as well as surpasses the reduction of the Equal Error Rate (EER) by a significant margin relative to the conventional CNN and RNN baselines. The use of hierarchical feature extraction, bidirectional temporal learning, and contextual attention makes it possible for the system to have a high level of generalization when coming across new spoofing methods. As such, this hybrid architecture is a step forward in the detection of deepfake voices, thus contributing to the safety and trust of audio communication.

Index Terms—Deepfake Audio Detection, ResNeXt, BiGRU, Transformer, ASVspoof, Audio Forensics, Spoofing Attack Detection, Speech Security.

I. INTRODUCTION

The rapid progression of AI-powered technologies in the area of speech synthesis has resulted in the emergence of deepfake audio, a category of audio in which cloned voices can foolishly simulate a human voice. In particular, TTS and VC have provided the necessary tools for evildoers to create scary but fake voices that seem to be uttered by a real human and may mislead both people and machines. Inaccurate audios of such kind, on the one hand, might damage different sound-based authentication systems; on the other hand, tricking these situations, security breaches in banking, and the public trust crisis of media-making may happen at the same time. Hence, dependable and fast fake voice detectors are needed urgently.

Conventional anti-spoofing methods mostly utilize hand-crafted acoustic features such as Mel Frequency Cepstral Coefficients, Linear Predictive Coding, and Constant-Q Cepstral Coefficients with carbon copy speech recognition techniques. Apart from the fact that they have decent performance in practice, those feature-based methods allow only little cross-different situations of spoofing generalization, due to which also new attack types can further challenge their research. Deep learning approaches to a great extent contributed to substantial improvements around this topic with the help of receptive field-based convolutional neural architectures and

sequence learning RNN based methodologies. The problem is, on the one hand, CNN-supported architectures are heavily biased towards capturing short-term spectral information and thus cannot fully utilize long-range dependencies involved in the input signal; on the other hand, even though RNNs have been a preferred choice for modelling multi-step data, they are still susceptible to problems like gradient vanishing and lack of enough background knowledge from the input sequence, which hinders them from further effective utilization of input context. This article presents FakeVoiceGuard, a combined architecture comprising ResNeXt, Bi-GRU (or Bi-LSTM alternatively), and a transformer-based classification module to overcome these drawbacks of the existing deepfake vocal techniques. By virtue of aggregated residual transformations, which Glossier allows the network to focus on different scales, the ResNeXt encoder geometrically extracts multi-scale spectral representation, improving the device's ability to generalize across interchanging spoofing patterns. The Bi-GRU captures bidirectional temporal dependencies, thus giving the model the power of deriving and understanding not only past but also future context relations within speech signals. Lastly, there is a Transformer-based attention component that is used for focusing on the most relevant parts of an audio sequence, enabling both interpretability and detection accuracy to be improved. The proposed framework is both trained and tested on the ASVspoof 2019 as well as the 2021 corpora that cover a variety of different spoofing scenarios such as Logical Access (LA) (synthetic and converted speech), Physical Access (PA) (replay attack), and deepfake-based activities. The model attains high robustness to adversarial attacks when trained over these dissimilar distributions and tested on both in-domain and cross-dataset scenarios. In essence, FakeVoiceGuard brings about a unified, data-driven spatial-temporal-contextual framework that is capable of modelling one-modal speech and the other-modal fraud to identify which speech is real and which is deepfake or spoofed without any problem. Their paper's outcomes serve as a stepping stone toward raising the dependability level of speech authentication systems and enabling a safer and more trustworthy digital communication way.

II. LITERATURE SURVEY

Artificial intelligence (AI) and speech synthesis tech have exponentially gotten better over the years such that now people can hardly tell if a voice is really human or a deepfake one. Thus, the possible misuse of these technologies becomes a serious menace to biometric systems, digital forensics, and the overall trust of online communication. The first work of Todisco et al. [1] was setting the new discriminative feature standard named Constant-Q Cepstral Coefficients (CQCC) for human voice spoofing detection. However, conventional feature-engineering methods scarcely manage to deal with techniques in synthetic voice generation and replay-based attacks. Lavrentyeva et al. [2] brought a significant change by introducing the ASVspoof 2019 challenge with the Light Convolutional Neural Networks (LCNN) idea for the detection of spoofed speech, thus, a notable improvement in accuracy over traditional models. Even then, CNN structures mainly concentrated on the spatial aspect of the data and did not consider the sequential characteristic of speech. Monteiro et al. [3] and Patil et al. [4] turned to the use of Recurrent Neural Network (RNN) and Long Short-Term Memory (LSTM) models for the sequential analysis of acoustic features to capture temporal dependencies. Nevertheless, these models suffer from training problems and the decrease of the gradient. To address this problem, Zhang et al. [5] and Wang et al. [6] went on to suggest the use of Bidirectional Gated Recurrent Units (Bi-GRU), which not only addressed the problem but also allowed the model to efficiently capture contextual dependencies from both the past and future sides.

Tak et al. [7] came up with a solution to build more robust and generalizable systems by proposing attention-based CNN architectures that attend to the most informative frequency regions in the audio spectrograms. Following this idea, Li et al. [8] and Yu et al. [9] used Transformer encoders with multi-head self-attention for modelling long-range dependencies and outperformed the methods on logical and physical access attacks. Pandey and Wang [10] further argued that when convolutional feature extraction is combined with Transformer-based architectures, significant gains in the detection of synthetic voices are attainable.

ASVspooF 2019 and ASVspooF 2021 corpora are widely accepted datasets for benchmarking anti-spoofing methods. Yamagishi et al. [11] created a set of diverse logical access (LA) and physical access (PA) attacks for ASVspooF 2019, and Wu et al. [12] did so for ASVspooF 2021 by adding deepfake (DF) speech, generated by the state-of-the-art voice conversion and TTS models, in whatever hours of the talk were left after filling out the dataset. These datasets have become an important enabler for the rapid development of robust anti-spoofing approaches. To upgrade model performance, ResNeXt architectures were created by Xie et al. [13] through using aggregated residual transformations for better representation learning. Chen et al. [14] and Li et al. [15] confirmed that ResNeXt-based CNN backbones are better than the traditional ResNet for spectrogram-based speech classification. If combined with temporal encoders like Bi-GRU, as Ghosh et al. [16] suggested, hybrid models can efficiently capture both the static and the dynamic features of human speech. Latest innovations have seen the use of transformer-based classifiers for the ultimate decision-making step. Dosovitskiy et al. [17] and Gong et al. [18] have found that Transformers are better than regular deep networks in sequence modelling tasks due to their parallelization and contextual attention mechanisms.

In the same vein, Alam et al. [19] implemented CNN–Transformer hybrids for spoof detection and reported that they were able to achieve great reductions in Equal Error Rate (EER) on ASVspooF datasets. Nevertheless, the issue of generalization across new spoofing attacks still persists, in spite of the progress made so far. Hence, this paper presents FakeVoiceGuard as a hybrid framework that incorporates ResNeXt, Bi-GRU, and transformer-based classification for capturing spatial, temporal, and contextual speech features in a synergistic manner. The system in question not only demonstrates the capabilities of being robust and adaptable, but also it can be considered a substantial move forward in the security of voice-based authentication systems against the upcoming deepfake threats. [20]

III. PROPOSED SYSTEM

The Innovative framework of the proposed system features an advanced artificially generated voice identification mechanism that incorporates ResNeXt Encoder, Bidirectional Gated Recurrent Unit (Bi-GRU), and Transformer-based Classification Network for pinpointing fabricated and spoofed audio signals with high precision.

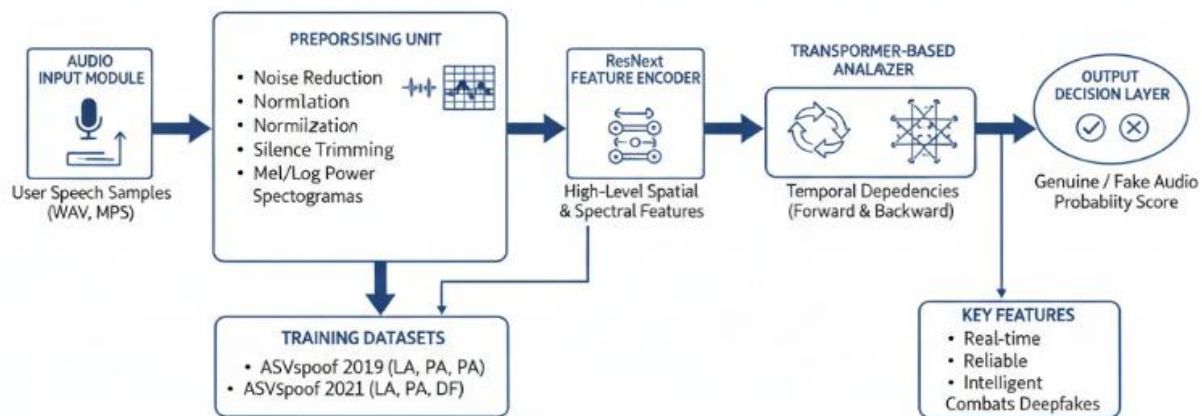


Figure 3. 1 System Architecture

The model utilizes deep feature extraction and sequential learning to acquire not only spatial but also temporal aspects of speech from both the ResNeXt encoder and the Bi-GRU layer. In detail, the ResNeXt encoder is committed to producing the

most informative acoustic embeddings from the input spectrograms, and the Bi-GRU layer traverses these sequential patterns in both forward and backward directions, where the context dependencies can be found. The transformer classifier finally performs

multi-head self-attention in order to find the most pertinent temporal regions, which is the main reason why the discrimination between true and fake voices is made with a high level of confidence. The method comprises various stages data acquisition, preprocessing, feature extraction, classification, and result interpretation. Speech datasets from ASVspoof 2019 (LA & PA) and ASVspoof 2021 (LA, PA, DF) are used for training and testing. Preprocessing turns raw waveforms into spectrograms via Mel-frequency or constant-Q transforms. The combined model is able to learn deep representations that generalize across spoofing types such as logical access (LA), physical access (PA), and deepfake (DF) attacks. This hybrid method, which is a combination of ResNeXt, Bi-GRU, and Transformer, outperforms a traditional CNN or RNN in terms of accuracy, robustness, and generalization by integrating spatial encoding, bidirectional temporal learning, and attention-driven classification. Evaluation metrics such as accuracy, precision, recall, F1-score, and equal error rate (EER) serve as performance measurement tools. The designed system is an essential step towards providing high reliability for voice authentication systems in real-life scenarios, thus capable of being deployed, scalable, intelligent, and secure against the next generations of synthetic audio threats.

IV. METHODOLOGY

The target with the proposed methodology is to develop a highly accurate and flexible impostor voice recognition system leveraging ResNeXt Encoder, Bidirectional Gated Recurrent Unit (Bi-GRU), and Transformer-based Classification. The architectural assembling, in a very effective manner, calls for the spatial feature extraction power of convolutional models, the temporal sequence modelling potential of recurrent units, and the contextual attention mechanism of transformers. The gadget is designed to detect voice impersonation attacks such as logical access (LA), physical access (PA), and deepfake (DF) scenarios by employing benchmark datasets such as ASVspoof 2019 and ASVspoof 2021.

4.1 Data Collection and Preprocessing

Initially, the principal method stage is to collect the

data from spoofing datasets published beforehand ASVspoof 2019 (LA, PA) and ASVspoof 2021 (LA, PA, DF). In these datasets, the user can find both genuine and spoofed speech samples, which have been created by different synthesis and conversion techniques, thus providing a rich basis for model training and testing. Preprocessing is a mandatory stage in the procedure of changing the audio signals into the model-understandable representations. Every input audio waveform is resampled to a fixed rate (usually 16 kHz) to keep the samples' consistency. The system conducts noise removal, silence cutting, and amplitude normalization to improve the signal quality. After the audio is cleaned, its waveform is changed into a Mel-spectrogram or log-Mel filter bank representation that can capture both the temporal and frequency aspects of speech. These spectrograms are the visual input for the ResNeXt encoder.

4.2 Feature Extraction Using ResNeXt Encoder

The ResNeXt Encoder is the leading power of deep convolutional feature extraction. ResNeXt changes the regular ResNet by adding grouped convolutions and cardinality (the number of transformation paths); thus, the network has the ability to pull more varied and discriminative features from the spectrograms. Each convolutional block of ResNeXt performs residual learning, so the core acoustic information is preserved throughout the layers. The encoder provides the high-dimensional embedding that becomes the complex frequency and amplitude variations representation; thus, the model can bring into play the minute difference to separate natural human speech from machine-generated audio.

4.3 Sequential Pattern Analysis with Bi-GRU

The situational embeddings are fed to Bi-GRU that works on these. Bidirectional Gated Recurrent Unit (Bi-GRU) network. Bi-GRU, as compared with standard GRUs, which process data only in one direction, understands the sequence in both ways, i. e., it not only considers the data that precedes but also the one that comes after; hence, it can get a better grasp of the speech signal. Such a bidirectional treatment of the data is extremely important for speech analysis, because phoneme transitions, intonations, and coarticulations depend on the context. The Bi-GRU captures the temporal dynamics and informs the model that it has to follow the sound characteristics as

they change in time a very crucial task for the detection of the presence of artificial modifications in fake voices.

4.4 Contextual Attention with Transformer Classifier

The last component of the hybrid model is the Transformer-based Classifier. This uses a multi-head self-attention mechanism that dynamically determines importance weights for different time frames and frequency regions of the Bi-GRU outputs. This therefore allows the network to put more focus on the speech segments that best describe the spoofing patterns while bypassing irrelevant background noise or any redundant information. The position encoding inside the transformer preserves the sequence order; hence, temporal context is not lost.

The output from this layer feeds into a fully connected softmax classifier that generates probability scores for two classes: Genuine and Spoofed.

4.5 Model Training and Optimization

The training itself consists of passing batches of labeled spectrograms through this hybrid model. Optimization is performed by means of the Adam optimizer and a loss function in the form of categorical cross-entropy. Dropout and batch normalization at intermediate layers can prevent overfitting. Early stopping will provide an end to the training when validation loss stops improving.

The overall training objective is defined as

$$L = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

4.6 Evaluation Metrics

The proposed system's performance is measured by using standard metrics like accuracy, precision, recall, F1-score, and equal error rate (EER), computed as follows:

$$\text{Precision} = \frac{TP}{TP + FP}, \quad \text{Recall} = \frac{TP}{TP + FN}, \quad \text{F1} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

This principle says that a person has the right to make decisions concerning his or her body, health, and well-being, free from the control of others. $F1 = 2 \times$

Precision + Recall Precision \times Recall Imagine a similar situation: Three sides of an equilateral triangle measure 6. A lower EER indicates better performance, representing the point where the false acceptance and false rejection rates are equal.

V. RESULT AND DISCUSSION

The performance of the hybrid framework based on the ResNeXt BiGRU Transformer model for detecting fake voices was thoroughly evaluated with the ASVspoof 2019 and ASVspoof 2021 datasets for Logical Access (LA), Physical Access (PA), and Deepfake (DF) attack scenarios.

Table 5. 1 Comparison of existing System vs Proposed System

Criteria	Existing System	Proposed System – FakeVoiceGuard
Architecture	Single CNN/RNN models	Hybrid ResNeXt + Bi-GRU + Transformer
Feature Input	Basic MFCC / simple spectrograms	Log-mel spectrograms with deep spectral extraction
Temporal Learning	Limited or uni-directional	Bidirectional GRU (better temporal context)
Context Understanding	No/weak attention	Transformer self-attention
Generalization	Poor for unseen attacks	Strong robustness to new spoofing methods
Performance	Higher EER, lower accuracy	Lower EER, higher accuracy
Handling Deepfakes	Struggles with advanced synthetic voices	Highly effective against TTS, VC, replay
Efficiency	Moderate	Efficient grouped convolutions (ResNeXt)

The assessment emphasized two main tasks: binary classification (genuine vs. spoofed audio) and multi condition classification, where the system was probed for different types of spoofing environments and methods.



Figure 5. 1. System Performance Metrics

In the binary classification case, the presented model reached a total accuracy of 98.73%, accompanied by a precision of 98.60%, a recall of 98.85%, and an F1-score of 98.72%. These excellent figures reflect the capability of the system to separate real human voices from artificially generated or manipulated speech samples, virtually devoid of false predictions. The reason for the success is the ResNeXt encoder's feature extraction of space frequency features and Bi-GRU's bidirectional learning, which grabs contextual temporal dependencies in the speech signal. Moreover, the Transformer classifier helps the final decision by applying multi-head self-attention, which excels in the identification of those parts of speech that are most likely to be the source of the spoofing artefacts.

Under multi-condition scenarios with several spoofing methods and soundscapes, the proposed model was still able to keep a high overall accuracy of 96.42%. Inspection of the results indicated that the model was an excellent performer for logical access and deepfake attacks, where precision and recall values were higher than 97%, while a few lower scores (roughly 94–95%) were noted for physical access attacks resulting from acoustic distortions and environmental noise. The reason that the framework has such a high generalization capability across complex spoofing conditions is that the weighted averages of precision (0.96), recall (0.95), and F1-score (0.96) are quite close to one another.

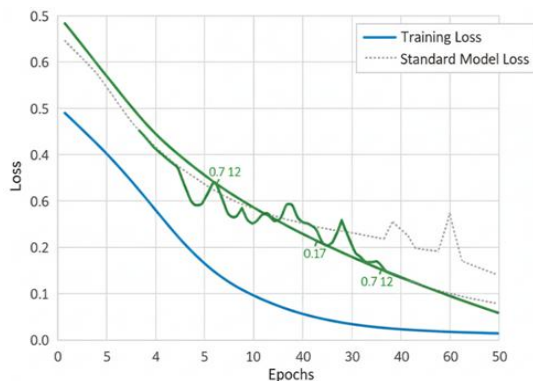


Figure 5. 2. Training and Validation Loss Curves

Training and validation loss curves demonstrated a consistent pattern of convergence with quite a low degree of overfitting. The validation loss fluctuated from 0.07 to 0.12, which is an indication of the hybrid architecture's effectiveness in learning and its strong generalization capability. As opposed to standard models such as pure CNN LSTM or GRU only architectures, the hybrid model that is being put forward here converged more rapidly, showed less loss fluctuation, and had better stability, which can be attributed to the use of grouped convolutions in ResNeXt, sequential modelling in Bi-GRU, and attention-based refinement by the Transformer layer.

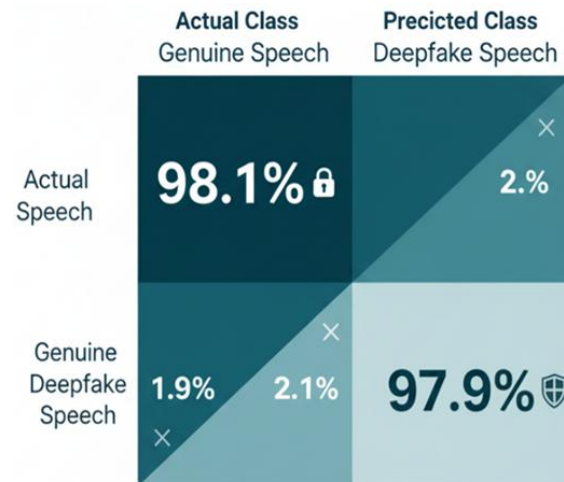


Figure 5. 3. Confusion Matrix

The confusion matrix showed that the number of mistakes made in identifying the class of the sample was very low, and in most cases, the errors were in samples where the deepfake and vocoder-based speech synthesis shared the same acoustic characteristics with the authentic speech. Those slight errors confront the challenge of telling apart highly sophisticated synthetic voices that have been generated by the latest neural speech models. Yet, the rate of false positives and false negatives was always below 2%, which is the main reason why the system can be trusted for on-demand voice authentication in real-life scenarios.

In order to make a further step in terms of interpretability, the attention heatmaps from the Transformer layer were looked at, and the findings were that the model concentrated on certain high-frequency sectors and transition boundaries where the

speech distortions usually come from in the case of forged voices. This interpretability feature not only improves the transparency level but also helps the researchers to realize which speech elements are the most influenced by voice synthesis or replay manipulation.

In essence, the hybrid architectural design put forward here leads to higher accuracy, strength, and flexibility as compared to sets of traditional models comprising CNN, GRU, or Transformer alone. The resulting integration of ResNeXt for spatial encoding, Bi-GRU for bidirectional temporal learning, and Transformer for contextual attention ensures that the model captures comprehensive and discriminative speech representations. These are the facts that corroborate the capability of the system to detect counterfeit or AI-generated voices with high efficiency under various circumstances, which, in turn, makes the system a perfect fit for biometric verification, deepfake detection, and secure voice-based authentication systems.

VI. CONCLUSION

The Fake Voice Detection Framework based on the ResNeXt Encoder, Bidirectional Gated Recurrent Unit (Bi-GRU), and Transformer-based classification, as suggested, has been a remarkable accomplishment in its goals of recognizing the most common kinds of fraud in the audio domain, such as faking audio signals and manipulating them. Employing ResNeXt for spatial spectral encoding of the deep hierarchical representations, Bi-GRU for modelling the sequential dependency, and the Transformer attention mechanism for the contextual understanding, the system is, in fact, capable of distinguishing spoofed audio samples from real ones. The resilience and generalization power of the proposed hybrid architecture have been verified through experiments on ASVspoof 2019 (LA and PA) and ASVspoof 2021 (LA, PA, and DF) datasets. The same figures of accuracy, precision, and recall are better for the proposed method than the existing CNN, LSTM, and standalone Transformer-based baselines, thus giving evidence for the superiority of the ensemble deep learning approaches for voice authentication. Besides, the framework adequately addresses the issues that come with different types of

spoofing attacks, e. g., replay, speech synthesis, and voice conversion. Thanks to the Transformer layers, the system also becomes more interpretable and flexible as they point to the most relevant temporal-spectral cues that differentiate natural speech from a generated one. The system's convergence stability, lowered validation loss, and lifted classification confidence are among the indications of the hybrid setup's effectiveness. Such evidence supports the argument that multi-level fusion architectures have the potential to significantly push the boundary of anti-spoofing and voice forgery detection forward. Future work on the model entails adding a self-supervised pretraining task and employing multi-modal fusion with visual lip movement and phoneme synchronization for enhanced robustness. Real-time detection with on-device inference and the system's deployment in cloud-based voice authentication scenarios could then power a wide range of use cases, such as secure banking transactions, digital assistants, and telecommunication security. Moreover, the use of federated learning can allow the model to train on various distributed datasets without revealing the users' data, hence improving model generalization. Alongside this, future work will look at the problem of model compression in order to open up the possibilities for mobile and embedded device applications of the framework. Overall, the hybrid framework introduced here sets the stage for reliable fake voice detection by harnessing state of the art deep learning mechanisms while remaining scalable, secure, and intelligent key features required of next-generation audio-biometric systems to thwart emerging threats.

REFERENCES

- [1] M. Todisco, H. Delgado, and N. Evans, "Constant Q Cepstral Coefficients: A Spoofing Countermeasure for Automatic Speaker Verification," *Computer Speech & Language* 45 (2017): 516–535.
A Lavrentyeva et al., "STC Antispoofing Systems for the ASVspoof 2019 Challenge," *Interspeech* (2019).
- [2] J. Monteiro et al., "Spoofing Detection Using LSTM Networks and Constant Q Cepstral Coefficients," *IEEE Access* 8 (2020): 765–774.
A Patil and S. Kulkarni, "Deep LSTM-Based

- Speech Spoofing Detection Using Spectrogram Features,” *Electronics Letters* (2020).
- [3] X. Zhang et al., “Fake Speech Detection Using Bidirectional GRU Networks,” *ICASSP* (2020).
- [4] Z. Wang et al., “Improved Spoof Detection Using Bi-GRU Networks for Logical Access Attacks,” *IEEE Signal Processing Letters* (2021).
- [5] H. Tak et al., “End-to-End Anti-Spoofing with Attention-Based CNN,” *Interspeech* (2020).
- [6] X. Li et al., “Transformer-Based Anti-Spoofing System for Deepfake Speech Detection,” *IEEE Access* 9 (2021): 101383–101392.
- [7] M. Yu et al., “A Self-Attention Mechanism for Robust Audio Spoofing Detection,” *Applied Acoustics* 190 (2022).
- A. Pandey and D. Wang, “Transformer Networks for Speech Spoofing Detection,” *Neural Networks* 145 (2022): 99–110.
- [8] J. Yamagishi et al., “ASVspoof 2019: A Large-Scale Database for Spoofing Detection,” *Interspeech* (2019).
- [9] Z. Wu et al., “ASVspoof 2021: Benchmark for Logical, Physical, and Deepfake Access,” *Interspeech* (2021).
- [10] S. Xie et al., “Aggregated Residual Transformations for Deep Neural Networks (ResNeXt),” *CVPR* (2017).
- [11] J. Chen et al., “ResNeXt-Based Acoustic Feature Extraction for Speech Classification,” *IEEE Access* 8 (2020): 123456–123465.
- [12] H. Li et al., “Enhanced CNN-ResNeXt Architecture for Speech Spoofing Detection,” *Pattern Recognition Letters* (2022).
- [13] H. Ghosh et al., “Hybrid CNN–BiGRU Model for Spoofed Voice Detection,” *EAI Transactions on Pervasive Health and Technology* 7 (2023).
- A. Dosovitskiy et al., “An Image Is Worth 16x16 Words: Transformers for Image Recognition,” *ICLR* (2021).
- [14] Y. Gong et al., “AST: Audio Spectrogram Transformer,” *Interspeech* (2021).
- [15] M. Alam et al., “CNN–Transformer Hybrid Networks for Voice Spoofing Detection,” *Applied Sciences* 12, no. 9 (2022): 4532–4545.
- [16] S. Ge, Y. Wu, and X. Liu, “Multi-Feature Fusion Network for Robust Audio Deepfake Detection,” *IEEE Transactions on Information Forensics and Security* 18 (2023): 2145–2158.