# IEEE Copyright Notice

arXiv:2512.21702v1 [cs.SD] 25 Dec 2025

# Zero-Shot to Zero-Lies: Detecting Bengali Deepfake Audio through Transfer Learning

Most. Sharmin Sultana Samu*, Md. Rakibul Islam†, Md. Zahid Hossain‡,
Md. Kamrozzaman Bhuiyan§, Farhad Uz Zaman¶
*Department of CSE, BRAC University, Bangladesh.
† ‡Department of CSE, Ahsanullah University of Science and Technology, Bangladesh.
§Enosis Solutions, Bangladesh.
¶Department of CSE, Southeast University, Bangladesh.
Email: *sharminsamu130@gmail.com, †rakib.aust41@gmail.com, ‡zahidd16@gmail.com,
§kamrozzamaan@gmail.com, ¶farhad.zaman@seu.edu.bd

*Abstract*—The rapid growth of speech synthesis and voice conversion systems has made deepfake audio a major security concern. Bengali deepfake detection remains largely unexplored. In this work, we study automatic detection of Bengali audio deepfakes using the BanglaFake dataset. We evaluate zero-shot inference with several pretrained models. These include Wav2Vec2-XLSR-53, Whisper, PANNsCNN14, WavLM and Audio Spectrogram Transformer. Zero-shot results show limited detection ability. The best model, Wav2Vec2-XLSR-53, achieves 53.80% accuracy, 56.60% AUC and 46.20% EER. We then fine-tune multiple architectures for Bengali deepfake detection. These include Wav2Vec2-Base, LCNN, LCNN-Attention, ResNet18, ViT-B16 and CNN-BiLSTM. Fine-tuned models show strong performance gains. ResNet18 achieves the highest accuracy of 79.17%, F1 score of 79.12%, AUC of 84.37% and EER of 24.35%. Experimental results confirm that fine-tuning significantly improves performance over zero-shot inference. This study provides the first systematic benchmark of Bengali deepfake audio detection. It highlights the effectiveness of fine-tuned deep learning models for this low-resource language.

*Keywords*—Bengali deepfake audio detection, audio forgery detection, synthetic audio detection, fake voice detection, Wav2Vec2

## I. INTRODUCTION

Audio deepfakes create false and misleading content that can deceive individuals or influence public opinion. They are also used in cybersecurity attacks that target commercial systems. Such audio forgeries are generated with text-to-speech (TTS), voice cloning (VC) or transformation and large language models (LLMs) for text-to-audio synthesis. These technologies make synthetic speech more natural and harder to detect.

Research on audio deepfake detection aims to design effective methods with limited labeled data. Public challenges [1]–[3] have driven much progress and formed an active research community. Early studies focused on time-frequency features [4] or specialized neural architectures [5], [6]. Recent works have applied large self-supervised models such as Wav2Vec2 [7] and WavLM [8], showing strong results on benchmark datasets [9]. However, challenge-based methods often lack industrial robustness because they rely on small datasets and perform poorly in open conditions. Newer competitions such as the Kaggle Deep Fake Detection Challenge and AISG Trusted Media Challenge [10] are closer to practical applications, but they emphasize video or multimodal data rather than audio alone. A further problem in this field is the open-set nature of detection. Both real and fake classes evolve, yet most studies assume a supervised setup. As new synthesis techniques emerge, a key challenge is to design systems that adapt quickly with minimal labeled data.

Most existing works in audio deepfake detection focus on English and a few high-resource languages. Studies on low-resource languages such as Bengali remain scarce. Bengali is spoken by more than 230 million people, yet resources for detecting synthetic audio in this language are limited. The lack of large annotated datasets and pretrained models tuned for Bengali makes the problem more challenging. This gap creates risks for social, financial and political misuse of deepfake audio in Bengali. To address this, we study Bengali audio deepfake detection using the publicly available BanglaFake [11] dataset. Our goal is to benchmark multiple architectures under zero-shot and fine-tuned settings.

The central research question of this work is how effectively existing deep learning models can detect Bengali audio deepfakes under both zero-shot and fine-tuned settings. We ask whether large pretrained models can generalize to Bengali without task-specific training and to what extent fine-tuning improves performance when trained on a domain-specific dataset. We also investigate which architectures, including convolutional, recurrent, residual and transformer-based models, are most suitable for Bengali deepfake detection.

Our key contributions mark the following:

- We present the first systematic study on Bengali audio deepfake detection using the BanglaFake dataset.
- We evaluate zero-shot inference with multiple pretrained models and analyze their limitations.
- We fine-tune six architectures, including CNN-based, ResNet, Transformer and hybrid models, for Bengali deepfake audio detection.
- We provide a comprehensive comparison of results using accuracy, precision, recall, F1 score, AUC and EER.
- We highlight the potential of fine-tuned deep learning models for building robust detection systems in low-

resource languages.

The paper is organized as follows. Section II reviews existing works on audio deepfake detection and related studies in speech processing. Section III provides background study relevant to this research. Section IV describes the proposed methodology, including zero-shot inference and fine-tuning strategies with different models. Section V presents the experimental results with detailed performance analysis. Section VI concludes the paper and highlights possible directions for future research in Bengali audio deepfake detection.

## II. RELATED WORKS

Many studies apply deep learning to audio deepfake detection. ResNet18 is widely used due to its simplicity and strong performance. It is used with feature engineering techniques such as LFCC and CQCC [12], [13]. Some models use attention mechanisms and recurrent layers to handle temporal dependencies. For example, a ResNet18-LSTM combination with multi-head attention improves robustness in noisy conditions [13]. RawNet2 and its variants are applied for end-to-end detection and feature extraction [14]–[16]. Transformer-based models like AASIST, ViT and AST show strong global context learning [17]–[19]. Self-supervised learning models such as Wav2Vec2, HuBERT, Whisper and AudioMAE provide good feature representations for downstream classifiers [9], [20]–[22]. Hybrid approaches using multimodal inputs also show promise. For instance, a combination of Wav2Vec2 for audio and mBERT for lyrics improves detection in musical deepfakes [17]. Some models use text-audio contrastive learning for zero-shot and multilingual scenarios [23].

Datasets play a critical role in model evaluation and generalization. ASVspoof 2019 and 2021 are the most commonly used datasets [12]–[14], [16], [18], [24], [25]. These datasets provide logical and physical access scenarios for controlled testing. Some studies use Fake or Real [25], [26] and WaveFake [15], [16] to evaluate performance on image-based and spectrogram-based inputs. Others focus on more diverse and real-world datasets. For example, in-the-wild corpora are introduced to assess performance under uncontrolled conditions [14], [21]. New datasets such as FakeMusicCaps [17] and SynHate [27] expand the scope to musical and hate speech detection. SynHate includes 37 languages, supporting multilingual evaluation. Other works use large-scale paired text and audio data for low-resource languages [23]. EVDA is used to test continual learning across eight deepfake datasets [28].

Performance varies depending on model design and input features. Temporal CNN on mel-spectrograms achieves 92% accuracy, outperforming traditional classifiers like SVM and Random Forest [26]. Whisper-small performs well on multilingual hate speech detection, reaching 85.4% accuracy [27]. SpecRNet achieves performance close to LCNN but with fewer parameters and faster inference [15]. Adaptive adversarial training reduces EER for LCNN and RawNet3 under white-box attacks [16]. RegO improves continual learning by

reducing forgetting and enhancing generalization [28]. Multi-view and multi-scale feature fusion improves detection in noisy and short utterance conditions [20], [21]. Despite these advances, several models show sharp performance drops in real-world or out-of-domain scenarios [14]–[16], [18]. Short audio clips, unseen attacks and transferability issues reduce model reliability.

Limitations are commonly reported across the reviewed papers. Many models rely on handcrafted features, which may not adapt to new spoofing methods [12], [24]. Some models are sensitive to speaker, language or recording conditions [18], [23], [27]. CNN-based models often lack robustness to adversarial examples [29]. Training and inference are often resource-intensive, especially for large transformers or ensemble models [18], [22]. Generalization to novel attacks or domains remains a major challenge [16], [28]. Few studies explore interpretability or explainability of detection results [25]. Human performance comparisons reveal shared weaknesses between AI systems and users, especially against TTS attacks [24]. Evaluation is often limited to specific datasets or attack types, reducing cross-study comparability.

Future work in this field points to several promising directions. Many studies propose stronger generalization through multi-task and continual learning [19], [28], [30]. Researchers suggest integrating diverse features, including raw audio, spectral inputs and multimodal information [17], [20], [21]. Improvements in multilingual performance are recommended, especially using self-supervised and transformer-based models [18], [23], [27]. Lightweight architectures with fewer parameters are also encouraged for real-time deployment [15], [22]. Adaptive detection mechanisms can help address evolving spoofing strategies [12], [16]. Advanced adversarial training and robust evaluation protocols are needed to improve security [16], [29]. Standardization of datasets and benchmarks is another critical step for consistent evaluation [18].

Furthermore, current benchmarks and datasets are limited in linguistic diversity. Notably, no prior work has explicitly focused on Bengali, a major low-resource language with millions of native speakers. This gap restricts the applicability of current systems in real-world multilingual contexts. To address these limitations, our research aims to develop a robust, Bengali-capable audio deepfake detection system with improved generalization, cross-lingual transferability and resilience to adversarial conditions.

## III. BACKGROUND STUDY

### A. Pretrained Models for Zero-Shot Audio Deepfake Detection

Wav2Vec2-XLSR-53 [31] is a self-supervised speech representation model that learns contextual audio embeddings from raw waveform data. Whisper-small and Whisper-medium [32] are end-to-end speech recognition models that convert audio to text while producing robust audio features useful for downstream tasks. PANNsCNN14 [33] is a convolutional neural network trained on audio spectrograms for sound event detection and generates feature embeddings for classification.

WavLM-Base-Plus [8] is a self-supervised speech model designed to capture both acoustic and linguistic information from raw audio. Audio Spectrogram Transformer (AST) [34] is a transformer-based audio spectrogram model pretrained on AudioSet, optimized to detect various sound events.

### B. Deep Learning Models for Fine-Tuned Bengali Audio Deepfake Detection

Wav2Vec2-Base [7] is a self-supervised speech representation model that learns contextual audio embeddings from raw waveform and can be adapted for classification tasks. LCNN is a light convolutional neural network designed to extract time-frequency features from spectrograms for audio classification. ResNet18 is a residual network that captures hierarchical audio features through skip connections and deep convolutional layers. LCNN-Attention extends LCNN by adding an attention mechanism to focus on important regions in the audio feature maps. ViT-B16 is a transformer-based model that divides spectrograms into patches and processes them with self-attention for audio representation learning. CNN-BiLSTM combines convolutional layers for local feature extraction and bidirectional LSTM layers for capturing temporal dependencies in audio sequences.

### C. Evaluation Metrics

Accuracy represents the percentage of correctly classified samples among all samples and indicates overall performance. Precision measures the proportion of true positive predictions among all positive predictions and reflects the reliability of positive detections. Recall calculates the proportion of true positives detected among all actual positive samples and indicates the model's sensitivity. F1 Score is the harmonic mean of precision and recall and balances both false positives and false negatives. Equal Error Rate (EER) is the point where false acceptance rate equals false rejection rate and is a critical metric for assessing system robustness in spoof detection. A smaller EER indicates better model performance, while a larger EER shows weaker discrimination. Area Under the Curve (AUC) measures the ability of the model to distinguish between classes across different thresholds and indicates the overall discrimination capability of the classifier. A larger AUC indicates stronger discrimination capability, while a smaller AUC shows poor separability between classes.

## IV. METHODOLOGY

### A. Dataset and Preprocessing

We have used the publicly available BanglaFake [11] audio dataset. The BanglaFake dataset is a benchmark resource for developing and evaluating Bengali deepfake audio detection models. It contains 12,260 real and 13,260 synthetic speech samples in WAV format, each lasting approximately 6–7 seconds and recorded at 22,050 Hz. Real audio is sourced from the SUST TTS Corpus [35] and Mozilla Common Voice [36], covering seven speakers, while deepfake audio is generated using a VITS-based text-to-speech model trained on the SUST TTS Corpus, designed to mimic human speech with subtle artifacts. The dataset follows the LJ Speech format which provides standardized metadata and naming conventions to ensure compatibility with existing TTS and audio processing tools. This organization supports robust training, evaluation and benchmarking of deepfake detection systems in Bengali. All audio files are resampled to 16 kHz to maintain consistency. Audio samples are truncated or zero-padded to a fixed duration depending on the model requirements. Mel-spectrograms are extracted with 64–128 Mel bands, using standard FFT and hop lengths. Spectrograms are converted to decibel scale and normalized. For image-based models, spectrograms are resized to 224×224 pixels and replicated across three channels. A custom PyTorch Dataset class handles batch loading and preprocessing.

### B. Zero-Shot Inference Models

We apply zero-shot inference using six pretrained models: Wav2Vec2-XLSR-53, Whisper-small, Whisper-medium, PANNsCNN14, WavLM-Base-Plus and Audio Spectrogram Transformer (AST). Each model is used without fine-tuning. Audio inputs are preprocessed as required for each model. The models generate embeddings or predictions which are evaluated directly using standard metrics. Figure 1 illustrates our proposed methodology.
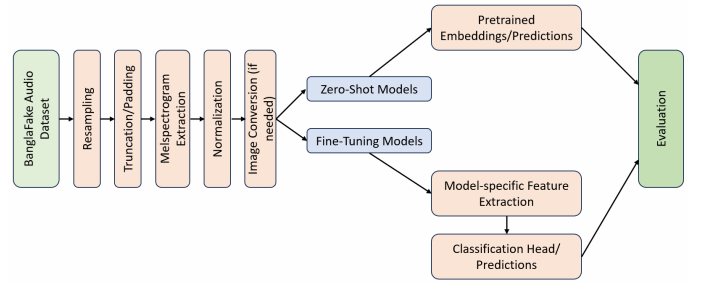


Fig. 1: Proposed Methodology

### C. Fine-Tuning Models

For fine-tuning, six models are trained on BanglaFake [11] dataset: Wav2Vec2-Base, LCNN, LCNN-Attention, ResNet18, ViT-B16 and CNN-BiLSTM.

- **Wav2Vec2-Base**: Pretrained Wav2Vec2 features are frozen. The classification head is trained on fixed-length 5-second audio samples. Outputs are binary logits for real or fake audio.
- **LCNN**: Mel-spectrograms are transformed into [1, 64, time] tensors. The network uses MFM layers and max pooling. A fully connected layer maps features to two classes.
- **LCNN-Attention**: Attention is applied over the time axis to highlight important temporal features. Outputs pass through fully connected layers for binary classification.
- **ResNet18**: Mel-spectrograms are expanded to 3-channel images. Features are extracted using a pretrained ResNet18 backbone. A classifier layer produces logits for deepfake detection.

- **ViT-B16**: Mel-spectrogram images are divided into patches for a Vision Transformer. Features are processed through the transformer encoder and the final classification head outputs binary logits.
- **CNN-BiLSTM**: Spectrograms are split into windows, converted to images and passed through a pretrained ResNet18 backbone. The resulting feature sequence goes through a bidirectional LSTM, followed by fully connected layers for classification.

Performance is measured using accuracy, precision, recall, F1-score, Equal Error Rate (EER) and Area Under the Curve (AUC).

### D. Experimental Setup

The dataset is split into training, validation and testing sets with a ratio of 70:15:15. Two experimental settings are designed: zero-shot inference using pretrained models and fine-tuning with the BanglaFake dataset. Training setups employ Adam-based optimization, weighted or standard cross-entropy losses and early stopping criteria. Model checkpoints are preserved based on validation performance.

TABLE I: Comparison of deepfake audio detection models.

| Model Name | Sequence Modeling | Audio Duration | Loss Function |
|---|---|---|---|
| LCNN | No | 5 sec | Cross-Entropy |
| LCNN-Attention | Yes | 5 sec | Cross-Entropy |
| ResNet18 | No | 5 sec | Binary Cross-Entropy |
| ViT-B16 | No | 5 sec | Cross-Entropy |
| CNN-BiLSTM | Yes | 5 sec | Weighted Binary Cross-Entropy |
| Wav2Vec2-Base | Yes | 5 sec | Cross-Entropy |

Table I compares several fine-tuned models used for Bengali deepfake audio detection. All models are trained with five-second audio inputs to maintain consistency. LCNN, ResNet18 and ViT-B16 rely on frame-level representations without sequence modeling, which limits their ability to capture temporal patterns in speech. LCNN-Attention, CNN-BiLSTM and Wav2Vec2-Base integrate sequence modeling, which helps them learn dependencies across time and improves detection of subtle manipulation cues. The choice of loss functions also reflects model design. Standard cross-entropy is applied in most cases as it is effective for classification tasks. ResNet18 uses binary cross-entropy, which simplifies the decision boundary for two-class detection. CNN-BiLSTM employs weighted binary cross-entropy to handle class imbalance, as fake and real audio data are often unevenly distributed.

TABLE II: Training hyperparameters for deepfake audio detection models.

| Model Name | Optimizer | Learning Rate | Batch Size | Epochs |
|---|---|---|---|---|
| LCNN | Adam | 0.0001 | 32 | 2 |
| LCNN-Attention | Adam | 0.0001 | 16 | 14 |
| ResNet18 | Adam | 0.0001 | 32 | 3 |
| ViT-B16 | Adam | 0.0001 | 8 | 3 |
| CNN-BiLSTM | Adam | 0.0001 | 8 | 10 |
| Wav2Vec2-Base | AdamW | 0.00005 | 4 | 1 |

Table II reports the training hyperparameters used for different models in Bengali deepfake audio detection. All models use the Adam optimizer except Wav2Vec2-Base, which adopts AdamW for better regularization. The learning rate remains fixed at 0.0001 for most models to ensure stable convergence, while Wav2Vec2-Base uses a lower rate of 0.00005 due to its large parameter size. Batch sizes vary across models and reflect the trade-off between memory usage and training stability. Larger batch sizes are applied to lightweight models such as LCNN and ResNet18, while smaller batches are required for ViT-B16, CNN-BiLSTM and Wav2Vec2-Base because of higher computational demands. The number of epochs also differs and indicates model complexity and convergence behavior. LCNN and ResNet18 converge quickly within few epochs, while LCNN-Attention and CNN-BiLSTM require longer training to capture temporal patterns. Wav2Vec2-Base is trained for only one epoch due to computation resource constraint.

## V. RESULT ANALYSIS

In this section, we present the experimental results of both zero-shot inference models and fine-tuned models.

### A. Zero-Shot Classification

The zero-shot models showed moderate to low performance. The highest accuracy was achieved by Wav2Vec2-XLSR-53 at 53.8%. PANNsCNN14 reached an accuracy of 50% but had perfect recall (100%), indicating it correctly identified all positive samples but also misclassified some negative samples. Whisper-small and Whisper-medium had lower accuracy (48.2% and 46%). Their precision and recall were equal, suggesting balanced performance but limited classification capability. WavLM-Base-Plus had zero precision and recall, indicating failure in correct positive predictions. The EER values ranged from 46.2% to 60.0%. The AUC values were generally low, with Wav2Vec2-XLSR-53 achieving the highest at 56.6%. Overall, zero-shot models showed limited ability to classify correctly without task-specific training. Table III reports zero-shot classification results.

TABLE III: Performance comparison of pre-trained audio models for zero-shot classification. Acc, Prec, Rec, F1, EER and AUC stands for Accuracy, Precision, Recall, F1 Score, Equal Error Rate and Area Under the Curve respectively. Values are presented in percentage.

| Model Name | Acc | Prec | Rec | F1 | EER | AUC |
|---|---|---|---|---|---|---|
| Wav2Vec2-XLSR-53 | 53.8 | 52.0 | 53.8 | 52.9 | 46.2 | 56.6 |
| Whisper-small | 48.2 | 48.2 | 48.2 | 48.2 | 51.8 | 47.7 |
| Whisper-medium | 46.0 | 46.0 | 46.0 | 46.0 | 54.0 | 44.5 |
| PANNsCNN14 | 50.0 | 50.0 | 1.0 | 66.7 | 52.6 | 43.2 |
| WavLM-Base-Plus | 50.0 | 0.0 | 0.0 | 0.0 | 60.0 | 36.4 |
| AST | 40.1 | 38.3 | 40.1 | 39.2 | 59.9 | 36.9 |

### B. Fine-Tuned Classification

Fine-tuned models demonstrated higher performance across all metrics. Table IV presents classification results of fine-tuned models.

TABLE IV: Performance comparison of deepfake audio detection fine-tuned models.

| Model Name | Acc | Prec | Rec | F1 | EER | AUC |
|---|---|---|---|---|---|---|
| Wav2Vec2-Base | 65.28 | 53.37 | 98.49 | 69.23 | 30.58 | 76.17 |
| LCNN | 48.36 | 46.73 | 52.53 | 49.46 | 61.23 | 50.48 |
| ResNet18 | 79.17 | 65.66 | 99.53 | 79.12 | 24.35 | 84.37 |
| LCNN-Attention | 78.43 | 64.94 | 99.11 | 78.47 | 23.01 | 88.48 |
| ViT-B16 | 78.65 | 65.41 | 97.97 | 78.45 | 22.26 | 86.63 |
| CNN-BiLSTM | 78.49 | 64.92 | 99.53 | 78.58 | 29.76 | 79.63 |

ResNet18 achieved the highest accuracy (79.17%) and a very high recall (99.53%). LCNN-Attention and ViT-B16 also performed well, with accuracy above 78% and AUC above 86%. Wav2Vec2-Base showed moderate performance with 65.28% accuracy and high recall (98.49%). LCNN had lower accuracy (48.36%) but moderate F1 score (49.46%). The EER values for fine-tuned models were much lower than zero-shot models, with ViT-B16 at 22.26% and LCNN-Attention at 23.01%, indicating better separation between classes. Overall, fine-tuning significantly improved classification performance and reliability compared to zero-shot approaches.
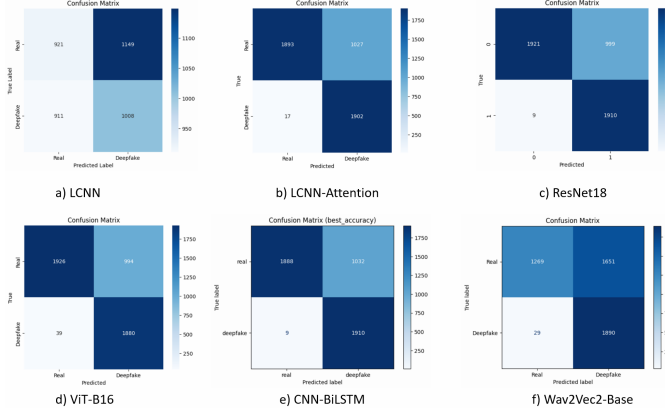


Fig. 2: Confusion matrix of the fine-tuned classification models.

Figure 2 presents the confusion matrix of six fine-tuned models. LCNN shows high misclassification with balanced errors in both classes. LCNN-Attention improves detection with very few errors for deepfake but higher confusion for real. ResNet18 achieves very strong performance with minimal misclassification in both classes. ViT-B16 also performs well with slightly higher misclassification for real but low error for deepfake. CNN-BiLSTM gives results close to ResNet18 with almost no misclassification for deepfake but higher confusion for real. Wav2Vec2-Base performs poorly for real detection with high misclassification but detects deepfake with strong accuracy.

Figure 3 presents the ROC curves of six fine-tuned models. LCNN shows poor separability with an AUC of 50.48% and the curve remains close to the diagonal. LCNN-Attention performs strongly with an AUC of 88.48% and the curve rises steeply toward the top left corner. ResNet18 achieves an AUC of 84.37% with a smooth curve showing consistent discrimi-
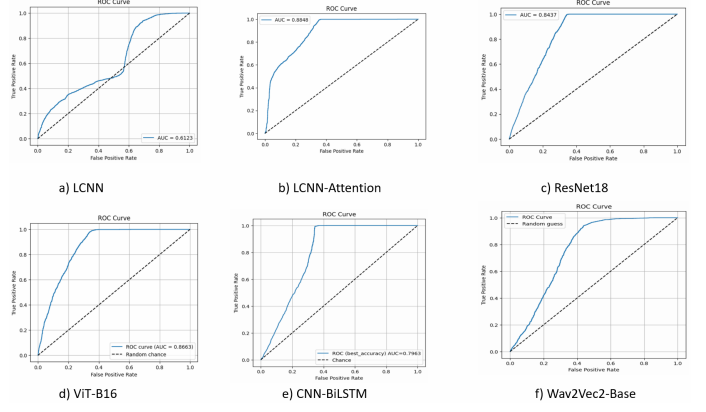


Fig. 3: ROC curve of the fine-tuned classification models.

nation ability. ViT-B16 provides one of the best results with an AUC of 86.63% and the curve remains close to the ideal boundary. CNN-BiLSTM produces an AUC of 79.63% with moderate classification strength and less steep rise. Wav2Vec2-Base shows weaker performance as the curve lies nearer to the diagonal with limited separation despite partial rise.
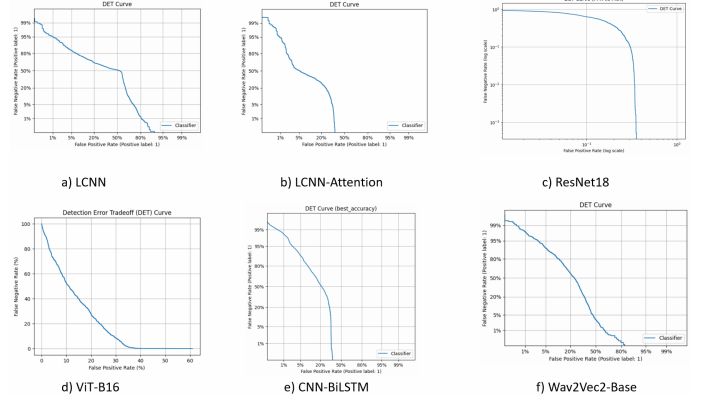


Fig. 4: DET curve of the fine-tuned classification models.

The Detection Error Tradeoff (DET) curves in Figure 4 show the error tradeoff across six models. LCNN presents high false negative rates across most thresholds, indicating weak separability. LCNN-Attention reduces errors with a smoother curve that drops faster, showing better balance between false positives and false negatives. ResNet18 achieves the strongest performance with a steep curve in logarithmic scale and very low error rates at optimal thresholds. ViT-B16 maintains a stable curve with gradual decline, reflecting reliable classification and low error at moderate thresholds. CNN-BiLSTM shows improved detection compared to LCNN, with a sharp drop in errors at lower false positive rates, but less steep than ResNet18. Wav2Vec2-Base produces higher error rates with a slower decline, indicating weaker discrimination power compared to vision-based models.

## VI. Conclusion and Future Works

We present the first systematic benchmark for Bengali deep-fake audio detection in this study. The results show that zero-shot inference with pretrained models provides limited effectiveness. Fine-tuned models achieve significant improvements across all metrics. ResNet18 gives the best performance among the tested architectures. The findings confirm that fine-tuning is necessary for robust detection in a low-resource setting. The study highlights the challenges of detecting deepfakes in Bengali audio and demonstrates the potential of deep learning methods for addressing this problem. Future work should expand the dataset to cover more speakers and diverse synthesis techniques. Cross-lingual transfer learning can be explored to improve performance in low-resource conditions. Robustness against adversarial attacks and unseen deepfake generation methods should be investigated. Lightweight models are required for real-time applications and deployment in resource-constrained environments. Future research should also integrate prosodic and linguistic features with acoustic cues to enhance detection reliability.

## References

[1] T. Kinnunen, M. Sahidullah, H. Delgado, M. Todisco, N. Evans, J. Yamagishi, and K. A. Lee, "The asvspoof 2017 challenge: Assessing the limits of replay spoofing attack detection," 2017.

[2] M. Todisco, X. Wang, V. Vestman, M. Sahidullah, H. Delgado, A. Nautsch, J. Yamagishi, N. Evans, T. Kinnunen, and K. A. Lee, "Asvspoof 2019: Future horizons in spoofed and fake audio detection," *arXiv preprint arXiv:1904.05441*, 2019.

[3] J. Yamagishi, X. Wang, M. Todisco, M. Sahidullah, J. Patino, A. Nautsch, X. Liu, K. A. Lee, T. Kinnunen, N. Evans, *et al.*, "Asvspoof 2021: accelerating progress in spoofed and deepfake speech detection," *arXiv preprint arXiv:2109.00537*, 2021.

[4] A. Fathan, J. Alam, and W. Kang, "Multiresolution decomposition analysis via wavelet transforms for audio deepfake detection," in *International Conference on Speech and Computer*, pp. 188–200, Springer, 2022.

[5] C.-I. Lai, N. Chen, J. Villalba, and N. Dehak, "Assert: Anti-spoofing with squeeze-excitation and residual networks," *arXiv preprint arXiv:1904.01120*, 2019.

[6] P. Aravind, U. Nechiyil, N. Paramparambath, *et al.*, "Audio spoofing verification using deep convolutional neural networks by transfer learning," *arXiv preprint arXiv:2008.03464*, 2020.

[7] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *Advances in neural information processing systems*, vol. 33, pp. 12449–12460, 2020.

[8] S. Chen, C. Wang, Z. Chen, Y. Wu, S. Liu, Z. Chen, J. Li, N. Kanda, T. Yoshioka, X. Xiao, *et al.*, "Wavlm: Large-scale self-supervised pre-training for full stack speech processing," *IEEE Journal of Selected Topics in Signal Processing*, vol. 16, no. 6, pp. 1505–1518, 2022.

[9] Y. Guo, H. Huang, X. Chen, H. Zhao, and Y. Wang, "Audio deepfake detection with self-supervised wavlm and multi-fusion attentive classifier," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 12702–12706, IEEE, 2024.

[10] W. Chen, S. L. B. Chua, S. Winkler, and S.-K. Ng, "Trusted media challenge dataset and user study," in *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*, pp. 3873–3877, 2022.

[11] I. A. Fahad, K. Asif, and S. Sikder, "Banglafake: Constructing and evaluating a specialized bengali deepfake audio dataset," *arXiv preprint arXiv:2505.10885*, 2025.

[12] T. Chen, A. Kumar, P. Nagarsheth, G. Sivaraman, and E. Khoury, "Generalization of audio deepfake detection.," in *Odyssey*, pp. 132–137, 2020.

[13] J. Pan, S. Nie, H. Zhang, S. He, K. Zhang, S. Liang, X. Zhang, and J. Tao, "Speaker recognition-assisted robust audio deepfake detection.," in *Interspeech*, pp. 4202–4206, 2022.

[14] N. M. Müller, P. Czempin, F. Dieckmann, A. Froghyar, and K. Böttinger, "Does audio deepfake detection generalize?," *arXiv preprint arXiv:2203.16263*, 2022.

[15] P. Kawa, M. Plata, and P. Syga, "Specrnet: Towards faster and more accessible audio deepfake detection," in *2022 IEEE International Conference on Trust, Security and Privacy in Computing and Communications (TrustCom)*, pp. 792–799, IEEE, 2022.

[16] P. Kawa, M. Plata, and P. Syga, "Defense against adversarial attacks on audio deepfake detection," *arXiv preprint arXiv:2212.14597*, 2022.

[17] Y. Li, Q. Sun, H. Li, L. Specia, and B. W. Schuller, "Detecting machine-generated music with explainability–a challenge and early benchmarks," *arXiv preprint arXiv:2412.13421*, 2024.

[18] J. Yi, C. Wang, J. Tao, X. Zhang, C. Y. Zhang, and Y. Zhao, "Audio deepfake detection: A survey," *arXiv preprint arXiv:2308.14970*, 2023.

[19] T. D. N. Le, K. K. Teh, and H. D. Tran, "Continuous learning of transformer-based audio deepfake detection," *arXiv preprint arXiv:2409.05924*, 2024.

[20] Y. Yang, H. Qin, H. Zhou, C. Wang, T. Guo, K. Han, and Y. Wang, "A robust audio deepfake detection system via multi-view feature," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 13131–13135, IEEE, 2024.

[21] Y. Chen, J. Yi, J. Xue, C. Wang, X. Zhang, S. Dong, S. Zeng, J. Tao, L. Zhao, and C. Fan, "Rawbmamba: End-to-end bidirectional state space model for audio deepfake detection," *arXiv preprint arXiv:2406.06086*, 2024.

[22] D. Combei, A. Stan, D. Oneata, and H. Cucu, "Wavlm model ensemble for audio deepfake detection," *arXiv preprint arXiv:2408.07414*, 2024.

[23] R. Ranjan, L. Ayinala, M. Vatsa, and R. Singh, "Multimodal zero-shot framework for deepfake hate speech detection in low-resource languages," *arXiv preprint arXiv:2506.08372*, 2025.

[24] N. M. Müller, K. Pizzi, and J. Williams, "Human perception of audio deepfakes," in *Proceedings of the 1st international workshop on deepfake detection for audio multimedia*, pp. 85–91, 2022.

[25] O. A. Shaaban, R. Yildirim, and A. A. Alguttar, "Audio deepfake approaches," *IEEE Access*, vol. 11, pp. 132652–132682, 2023.

[26] J. Khochare, C. Joshi, B. Yenarkar, S. Suratkar, and F. Kazi, "A deep learning framework for audio deepfake detection," *Arabian Journal for Science and Engineering*, vol. 47, no. 3, pp. 3447–3458, 2022.

[27] R. Ranjan, K. Pipariya, M. Vatsa, and R. Singh, "Synhate: Detecting hate speech in synthetic deepfake audio," *arXiv preprint arXiv:2506.06772*, 2025.

[28] Y. Chen, J. Yi, C. Fan, J. Tao, Y. Ren, S. Zeng, C. Y. Zhang, X. Yan, H. Gu, J. Xue, *et al.*, "Region-based optimization in continual learning for audio deepfake detection," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 39, pp. 23651–23659, 2025.

[29] M. Rabhi, S. Bakiras, and R. Di Pietro, "Audio-deepfake detection: Adversarial attacks and countermeasures," *Expert Systems with Applications*, vol. 250, p. 123941, 2024.

[30] X. Zhang, J. Yi, C. Wang, C. Y. Zhang, S. Zeng, and J. Tao, "What to remember: Self-adaptive continual learning for audio deepfake detection," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, pp. 19569–19577, 2024.

[31] A. Conneau, A. Baevski, R. Collobert, A. Mohamed, and M. Auli, "Unsupervised cross-lingual representation learning for speech recognition," *arXiv preprint arXiv:2006.13979*, 2020.

[32] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, "Robust speech recognition via large-scale weak supervision," in *International conference on machine learning*, pp. 28492–28518, PMLR, 2023.

[33] Q. Kong, Y. Cao, T. Iqbal, Y. Wang, W. Wang, and M. D. Plumbley, "Panns: Large-scale pretrained audio neural networks for audio pattern recognition," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 2880–2894, 2020.

[34] Y. Gong, Y.-A. Chung, and J. Glass, "Ast: Audio spectrogram transformer," *arXiv preprint arXiv:2104.01778*, 2021.

[35] A. Ahmad, M. R. Selim, M. Z. Iqbal, and M. S. Rahman, "Sust tts corpus: A phonetically-balanced corpus for bangla text-to-speech synthesis," *Acoustical Science and Technology*, vol. 42, no. 6, pp. 326–332, 2021.

[36] R. Ardila, M. Branson, K. Davis, M. Henretty, M. Kohler, J. Meyer, R. Morais, L. Saunders, F. M. Tyers, and G. Weber, "Common voice: A massively-multilingual speech corpus," *arXiv preprint arXiv:1912.06670*, 2019.