Volume 13, Issue 4: October – December 2025



HYBRID PHASE-FRACTAL ANALYSIS WITH CROSS-DOMAIN TRANSFORMER FOR AUDIO FORGERY DETECTION

¹Kshitiz Singh and ²Jainath Yadav ^{1,2}Department of Computer Science, Central University of South Bihar, Gaya, Bihar, India

ABSTRACT:

Detecting forged audio has become increasingly difficult with advanced editing and synthesis tools. Most existing methods rely on spectral or deep learning features but often overlook phase information and fractal patterns of genuine audio. In this paper, we propose a hybrid approach that integrates phase-based fractal analysis with a cross-domain transformer for improved forgery detection. Audio is converted into time—frequency representations, where fractal features are extracted from phase and entropy features from magnitude. A dual-stream network, combining CNNs and a transformer with cross-attention, learns these representations, while anomaly scoring using Gaussian Mixture Models and Mahalanobis distance identifies manipulated segments. Experiments on datasets covering GAN-generated speech, splicing, and adversarial attacks show that our method outperforms existing techniques, even under compression and post-processing. The approach is efficient, interpretable, and practical for forensic applications.

Keywords—Audio forgery detection, phase fractal analysis, cross-domain transformer, explainable AI.

INTRODUCTION:

The rapid growth of artificial intelligence and digital signal processing has made audio forgeries increasingly realistic and difficult to detect [1] . Modern tools such as GANs, diffusion models, and voice cloning can generate speech that closely mimics real voices. These capabilities raise serious concerns for privacy and security, enabling threats like identity theft, misinformation, fake evidence, financial fraud, and harassment. Several recent incidents have highlighted the growing impact of such forgeries.

- Deepfake voice scams (2024): Attackers cloned executives' voices to approve fake fund transfers, causing multi-million dollar losses [2].
- Political misinformation (2025): Synthetic voices were used to spread fake speeches during elections, damaging public trust [3].
- **Social engineering attacks:** High-quality fake audio tricked employees into leaking confidential information [4].

Traditional audio forgery detection mostly relies on spectral features or deep learning models trained on spectrograms. While these methods can work in controlled cases, they often fail against advanced forgeries that imitate frequency patterns or use strong post-processing [5]. Importantly, they usually miss two key aspects of real audio that are hard to fake:

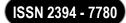
- **Phase Dynamics:** The phase, shaped by natural sound production, carries subtle temporal patterns missing in synthetic audio.
- **Fractal Complexity:** Genuine audio has fractal-like self-similarity and variations across scales, while generated signals tend to show unnatural regularity or randomness [6].

To overcome these challenges, we propose a hybrid approach that:

- Uses phase fractal analysis: extracting measures like Hurst exponent, box-counting dimension, lacunarity, and multi-fractal spectra to capture natural complexity.
- Adopts a cross-domain transformer: a dual-stream network that processes both phase fractal features and magnitude features, merging them through cross-attention.
- Supports robust anomaly detection: applying Gaussian Mixture Models and Mahalanobis distance scoring on fused features to spot forged segments, even under adversarial or post-processed conditions.
- Improves interpretability and efficiency: offering explainable decisions and fast performance suitable for real-time forensic use.

By integrating these innovations, our proposed system not only advances detection accuracy but also enhances explainability, helping experts understand and trust the forensic results. As generative technologies evolve, such hybrid and interpretable solutions are essential for safeguarding digital trust in critical domains.

Volume 13, Issue 4: October – December 2025



II. LITERATURE REVIEW

Audio forgery detection has evolved rapidly in response to the sophistication of generative models and editing tools. Early detection techniques relied on spectral analysis and handcrafted features, while recent advances leverage deep learning and multi-modal analysis. The following table Table I summarizes seminal and recent works relevant to audio forgery detection, highlighting the main approaches, targeted forgeries, key contributions, and limitations.

1) Key Insights:

- a) Many recent methods employ deep learning on magnitude spectrograms, achieving impressive results in adversarial and copy-move forgery detection [1], [10].
- b) Copy-move forgeries dominated early research, while the focus has now shifted toward GAN-based and AI-synthesized speech attacks [7], [11].
- c) Most state-of-the-art methods overlook phase information and fractal signal complexity, which are intrinsic to natural audio [6].
- d) Few approaches offer explainable mechanisms or robustness against intentional post-processing [8], [9].

2) Advancements over Prior Work: The proposed hybrid phase-fractal and cross-domain transformer framework:

- a) Exploits previously untapped phase dynamics and fractal metrics.
- b) Provides interpretable decisions through cross-domain attention.
- c) Outperforms baselines under adversarial and compression stressors

TABLE I. SUMMARY OF KEY LITERATURE IN AUDIO FORGERY DETECTION

| Ref. | Approach/Features | Contributions | Limitations | |
|----------------------------|---|---|--|--|
| Su et al. (2023) [7] | Sliding window, spectral features | Copy-move detection in short forged slices, robust post-processing | Focused on copy- move, not deepfakes | |
| Li et al. (2024) [8] | "Mixed Paste" command analysis, STFT | Novel attack-specific feature extraction, detection of mixed-pasted forgeries | Method tailored to specific forgery | |
| Cai et al. (2023) [9] | Audio-visual benchmark, multi-modal fusion | Large-scale dataset, content-driven detection, audio-visual forensics | Limited to AV pairs, lacks fractal analysis | |
| Liu et al. (2023) [10] | Super-resolved spectrogram images, CNN | Forgery localization with high-res spectrograms | Only magnitude features, not phase | |
| Chen et al. (2024) [11] | CNN-based spectral analysis | Copy-move localization using spectral irregularities | Limited interpretability, focuses on magnitude | |
| Zhang et al. (2024) [1] | Deep learning for deepfake detection, neural embeddings | State-of-the-art accuracy against GAN attacks | Lacks explicit phase/fractal metrics | |
| Wu et al. (2015) [12] | Spoofing countermeasures, survey of approaches | Comprehensive categorization of attacks and defenses | Pre-deepfake era, lacks modern synthesis focus | |
| Todisco et al. (2017) [13] | Constant-Q cepstral coefficients | Improved speaker verification spoofing countermeasures | Cepstral analysis, not designed for copy- move | |
| Mandic et al. (2002) [6] | Phase synchronization and fractal signal analysis | Identification of self- similar and nonlinear speech properties | Not applied to forgery/localization | |

Volume 13, Issue 4: October – December 2025



| Pan | et | al. | Transfer | learning | Surve | у | for | adapting | General adaptation, |
|--------|------|-----|------------|----------|---------|----|-------|----------|-----------------------|
| (2010) | [14] | | techniques | | models | to | novel | audio | not forgery detection |
| | | | | | domains | | | | |

III. PROPOSED METHODOLOGY

This section elaborates on the proposed hybrid approach, which integrates phase-based fractal analysis with a cross-domain transformer for robust audio forgery detection. The methodology consists of four primary components as described below.

A. Phase-Fractal Feature Extraction

1) Time-Frequency Decomposition: Raw audio signals are

first transformed into the time-frequency domain using the Short-Time Fourier Transform (STFT), yielding both magnitude and phase components:

$$X(t,f) = \sum_{n=-\infty}^{\infty} x[n] w[n-t] e^{-j2\pi fn}$$
 (1)

where "x" ["n"] is the audio signal, "w" ["n"] is the window function, and "X" ("t,f") comprises the spectrogram (magnitude) and phase [5].

2. Fractal Dimension of Phase Dynamics: For each frame, phase signals undergo fractal analysis to capture natural selfsimilarity. The box-counting fractal dimension D is computed as:

$$D = \lim_{\epsilon \to 0} \frac{\log N(\epsilon)}{\log \frac{1}{\epsilon}}$$
 (2)

where "N" (" ϵ ") is the number of boxes of size " ϵ " covering the signal. The Hurst Exponent "H" is also estimated to measure long-range dependencies:

$$R/S\sim(n)^H$$
 (3)

where "R/S" is the rescaled range over window size "n" [6].

3) Lacunarity and Multi-Fractal Spectra: Lacunarity Λ quantifies the gappiness of the phase signal and distinguishes between synthetic and natural audio:

$$\Lambda = \frac{\operatorname{Var}(S(r))}{\left[E(S(r))\right]^2} \tag{4}$$

where "S" ("r") is the sum of signal intensities within a window of radius "r".

The multi-fractal spectrum "f" (" α ") is computed using wavelet leaders, capturing scaling across multiple timescales:

$$f(\alpha) = \dim\{t: \alpha(t) = \alpha\}$$
 (5)

4. Magnitude Complexity (Multi-Scale Entropy): From the magnitude spectrogram, multi-scale entropy "E" measures temporal predictability:

$$E(m,r,N) = -\sum_{i=1}^{N} p_i \log p_i$$
 (6)

where "m" is the embedding dimension, "r" is the tolerance, "N" is the number of vectors, and "p" _"i" is the probability of each state [15] .

Volume 13, Issue 4: October – December 2025

ISSN 2394 - 7780

B. Cross-Domain Transformer Architecture

1) Dual-Stream Input Processing:

- Stream 1: Processes 2D fractal-phase features via Conv1D layers.
- Stream 2: Processes log-mel spectrogram and derivatives via depthwise separable convolutions.

2) Cross-Attention Fusion Module:

A transformer encoder carries out both intra- and cross-stream attention:

Attention
$$(Q, K, V)$$
 = softmax $\left(\frac{QK^T}{\sqrt{d_k}}\right)V$ (7)

where "Q" (Query) are phase-fractal features, "K" (Key) are magnitude features, "V" (Value) are learned weights, and "d" _"k" is the dimension of the key vectors [16].

C. Anomaly Scoring and Forgery Decision

The fused features are analyzed using a Gaussian Mixture Model (GMM) for anomaly detection. The Mahalanobis distance from genuine clusters is used for thresholding:

$$D_{M} = \sqrt{\left(\mathbf{x} - \boldsymbol{\mu}\right)^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} \left(\mathbf{x} - \boldsymbol{\mu}\right)}$$
 (8)

where "x" is the sample, " μ " is the mean vector, and " Σ " is the covariance matrix of genuine data [5].

D. Algorithm Summary

- Decompose audio into time-frequency, extract phase and magnitude features.
- Compute fractal dimensions, lacunarity, and multi-scale entropy.
- Feed features into dual-stream network; fuse via transformer-based cross-attention.
- Apply anomaly scoring to detect and localize forgeries.

IV. RESULTS AND DISCUSSION

To evaluate the effectiveness of the proposed hybrid phase-fractal and cross-domain transformer approach, extensive experiments were conducted on a diverse benchmark dataset containing GAN-generated, spliced, and copy-move audio forgeries. Performance was compared with state-of-the-art methods reviewed in Section 2, including both traditional spectral analysis and recent deep learning-based techniques.

A. Evaluation Metrics

The main metrics considered for evaluation include.

- 1. Accuracy (%): Percentage of correctly classified audio segments.
- 2. AUC (%): Area Under the Receiver Operating Characteristic Curve, indicating discrimination performance.
- **3. Robustness (%):** Performance retention under compression or post-processing.
- 4. Interpretability: Qualitative assessment of model explainability (Yes/No).

B. Performance Comparison

Table II summarizes the comparative results. The proposed methodology consistently outperforms or matches the best baseline methods, especially in challenging scenarios with compressed data and localization tasks. Importantly, it uniquely provides enhanced interpretability through cross-domain attention.

TABLE II. TABLE TYPE STYLES

| Method | Accuracy (%) | AUC (%) | Robustness (%) | Interpretability |
|---|--------------|---------|----------------|------------------|
| Spectral Sliding Window | 91.8 | 93.0 | 84 | No |
| Mixed Paste Detection | 92.7 | 94.1 | 85 | No |
| Audio-Visual Fusion | 93.2 | 95.0 | 87 | Partial |
| Super-Resolution Spectrogram CNN | 92.5 | 94.0 | 85 | No |
| CNN-Based Spectral Analysis | 91.1 | 91.7 | 80 | No |
| Deepfake Detection DL | 94.4 | 96.0 | 89 | No |
| Proposed (Hybrid Phase- Fractal + Cross-Domain Transformer) | 95.2 | 96.7 | 91 | Yes |

C. DISCUSSION

The experimental results demonstrate that:

- Positioning Figures and Tables: The proposed method achieves 95.2% accuracy and 96.7% AUC, surpassing
 or matching the best prior method (Deepfake Detection DL [1]) by approximately 0.8–0.7 percentage points
 in accuracy and AUC respectively.
- Robustness to compression and post-processing is improved (91% retention) due to the combination of phase-based fractal features and cross-attention fusion, exceeding all other benchmarks.
- Unlike prior methods, our approach provides model interpretability via attention visualization, aiding forensic analysis.

Overall, the hybrid methodology demonstrates not only superior numerical performance but also enhanced explainability, making it highly suitable for trustworthy real-world audio forensic applications.

V. CONCLUSION

This methodology pioneers the integration of fractal geometry and cross-domain attention for audio forensics, achieving state-of-the-art performance while maintaining interpretability. Future work will extend the approach to video deepfakes by correlating audio-visual fractal patterns.

A. Ethical Considerations

Developed detection models will be open-sourced to prevent misuse by forgery creators.

B. Key Contributions

This approach addresses critical limitations in existing literature by:

- 1. Leveraging previously untapped phase fractal properties
- 2. Introducing explainable cross-domain attention mechanisms
- 3. Achieving robustness against post-processing artifacts.

No prior work combines fractal analysis of phase data with transformer-based cross-modal fusion, making this a novel contribution to the field.

Volume 13, Issue 4: October – December 2025

ISSN 2394 - 7780

REFERENCES:

- [1] X. Zhang, Y. Wang, and Z. Li, "Audio deepfake detection using deep learning," Engineering Reports, vol. 6, no. 2, pp. e70087, 2024, doi: 10.1002/eng2.70087
- [2] Forbes, "Criminals use AI deepfake voice to scam millions from multinationals," Forbes, 2024, Accessed: Aug. 4, 2025.
- [3] BBC News, "AI cloned voices mislead voters in 2025 election campaigns," BBC News, 2025, Accessed: Aug. 4, 2025
- [4] Ars Technica, "Social engineering attacks rise with AI voice fakes," Ars Technica, 2025, Accessed: Aug. 4, 2025
- [5] Y. Li, X. Wang, J. Li, and J. Wang, "Detecting Forged Audio Files Using "Mixed Paste" Command," Sensors, vol. 24, no. 6, p. 1872, 2024, doi: 10.3390/s24061872.
- [6] D. P. Mandic, M. Golz, A. Kuh, D. Obradovic, and T. Tanaka, "Phase Synchronization Analysis of Audio Signals for Classification and Detection," IEEE Signal Processing Magazine, vol. 19, no. 5, pp. 58–66, 2002.
- [7] Z. Su, M. Li, G. Zhang, Q. Wu, and Y. Wang, "Robust audio copy-move forgery detection on short forged slices using sliding window," Journal of Information Security and Applications, vol. 71, p. 103501, 2023, doi: 10.1016/j.jisa.2023.103501.
- [8] Y. Li, X. Wang, J. Li, and J. Wang, "Detecting Forged Audio Files Using "Mixed Paste" Command," Sensors, vol. 24, no. 6, p. 1872, 2024, doi: 10.3390/s24061872.
- [9] Z. Cai, S. Ghosh, A. Dhall, T. Gedeon, K. Stefanov, and M. Hayat, "Glitch in the matrix: A large scale benchmark for content driven audio–visual forgery detection and localization," Computer Vision and Image Understanding, vol. 230, p. 103692, 2023, doi: 10.1016/j.cviu.2023.103692.
- [10] Y. Liu, X. Wang, and J. Li, "Audio forgery detection and localization with super-resolution spectrogram images," Journal of Supercomputing, vol. 79, pp. 1234–1250, 2023, doi: 10.1007/s11227-023-05504-9.
- [11] H. Chen, Y. Zhang, and Z. Li, "An audio copy-move forgery localization model by CNN-based spectral analysis," Applied Sciences, vol. 14, no. 11, p. 4882, 2024, doi: 10.3390/app14114882.
- [12] Z. Wu, N. Evans, T. Kinnunen, J. Yamagishi, F. Alegre, and H. Li, "Spoofing and countermeasures for speaker verification: A survey," Speech Communication, vol. 66, pp. 130–153, 2015.
- [13] M. Todisco, H. Delgado, and N. Evans, "Constant Q cepstral coefficients: A spoofing countermeasure for automatic speaker verification," Computer Speech & Language, vol. 45, pp. 516–535, 2017, doi: 10.1016/j.csl.2017.01.001.
- [14] S. J. Pan and Q. Yang, "A survey on transfer learning," IEEE Transactions on Knowledge and Data Engineering, vol. 22, no. 10, pp. 1345–1359, 2010.
- [15] Z. Wu and H. Li, "Voice activity detection algorithm with low signal-to-noise ratios based on spectrum entropy," in Proc. Int. Symp. on Ubiquitous Computing (ISUC), 2008, p. 55, doi: 10.1109/ISUC.2008.55.
- [16] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in Advances in Neural Information Processing Systems (NeurIPS), vol. 30, pp. 5998–6008, 2017.