

Speech Deepfakes Detection with Fast Fourier Transform using Complex Linear Algebra

Kevin Wirya Valerian - 13524019

Program Studi Teknik Informatika

Sekolah Teknik Elektro dan Informatika

Institut Teknologi Bandung, Jl. Ganesha 10 Bandung 40132, Indonesia

kevin.wirya.valerian@gmail.com 13524019@std.stei.itb.ac.id

Abstract—Nowadays, the growing realism of speech deepfakes demands detection methods grounded in fundamental signal theory rather than purely data-driven models. A well-known speech deepfake detection approach is based on the Fast Fourier Transform (FFT), rooted in complex linear algebra and geometric signal representations. Speech signals are modeled as vectors in complex vector spaces, where the FFT acts as a structured linear transformation projecting time-domain data onto orthonormal bases of complex exponentials. From a geometric perspective, genuine and synthetic speech can be represented differently in high-dimensional spectral spaces. Simple inconsistencies in phase behavior and energy distribution can be analyzed through algebraic measures derived from the complex FFT representation. This study demonstrates the applicability of FFT-based complex spectral analysis for distinguishing AI-generated and genuine speech.

Index Terms—Fast fourier transform, deepfakes detection, complex linear algebra, spectral geometry

I. INTRODUCTION

Recent advances in speech synthesis and voice conversion have enabled machines to generate highly realistic human speech. Modern text-to-speech and voice cloning systems are capable of producing audio that is increasingly difficult to distinguish from genuine human recordings. A study reports that human participants were only able to accurately distinguish real from AI-generated voices with an accuracy of 70.4%. This leads to a serious risks that can be posed by speech deepfakes in areas such as biometric authentication and digital forensics. Reports indicate that voice-based fraud has increased significantly in recent years, with financial losses reaching billions of dollars globally. Deepfake-related losses have already reached \$1.56 billion, with over \$1 billion occurring in 2025 alone.

While deep learning models have achieved impressive performance in generating natural-sounding speech, detecting such synthetic audio remains a challenging task. Many existing detection methods rely heavily on large neural networks trained on specific datasets. Although effective, these approaches often lack interpretability and tend to degrade when exposed to unseen deepfake generation methods. Moreover, purely data-driven models can be computationally expensive and difficult to analyze, making them less suitable and reliable at present time.

Speech signals, however, are fundamentally mathematical objects. A digital audio signal can be represented as a finite

sequence of samples, which naturally forms a vector in a high-dimensional space. Transforming this signal into the frequency domain using the Fast Fourier Transform (FFT) reveals its spectral structure, including phase behavior. These properties are governed by well-established principles of linear algebra and geometry, such as vector spaces, orthonormal bases, inner or dot products, and unitary transformations. From a linear algebra perspective, the FFT can be interpreted as a linear transformation that maps a time-domain signal into a complex frequency-domain vector. While synthetic speech models often reproduce spectral magnitudes effectively, differences in phase behavior may still arise.

Motivated by this observation, this paper proposes a speech deepfake detection approach grounded in the linear algebraic and geometric interpretation of the FFT. Rather than relying on large data-driven models, the proposed method focuses on analyzing basic algebraic properties of complex spectral representations. The approach aims to provide an interpretable demonstration of how concepts from complex linear algebra and geometry can be applied to distinguish genuine and synthetic speech signals.

II. THEORETICAL FRAMEWORK

A. Complex Numbers and Geometric Representation

A complex number is defined as

$$z = a + jb, \quad a, b \in \mathbb{R} \quad (1)$$

where a is the real part and b is the imaginary part, with $j = \sqrt{-1}$. Both i and j represent the same imaginary unit. We use j since it is commonly used for signal processing.

Complex numbers admit a geometric interpretation in the complex plane, where the real part corresponds to the horizontal axis and the imaginary part to the vertical axis. This representation allows us to visualize complex numbers as vectors emanating from the origin.

Every non-zero complex number can be expressed in polar form

$$z = re^{j\theta} \quad (2)$$

where

- $r = |z| = \sqrt{a^2 + b^2}$ is the **magnitude**, representing the distance from the origin

- $\theta = \arg(z) = \arctan(b/a)$ is the **phase**, representing the angle from the positive real (horizontal) axis

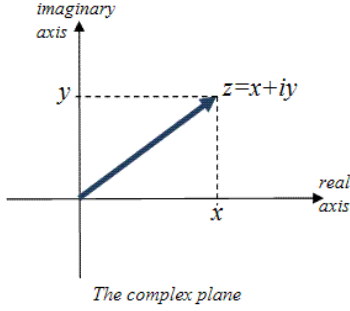


Fig. 1. Argand Plane (Complex Plane)
(Source: <https://helpingwithmath.com/complex-plane/>)

The exponential form relates to trigonometry via Euler's formula.

$$e^{j\theta} = \cos \theta + j \sin \theta \quad (3)$$

The FFT spectrum of audio signals consists of complex-valued coefficients. Each coefficient X_k can be decomposed into magnitude and phase components. The magnitude $|X_k|$ represents the energy at frequency bin k , while the phase $\angle X_k$ encodes temporal structure. A frequency bin is a specific range on the frequency axis used to group and analyze data. Our detection method exploits the observation that human speech exhibits **regular** phase structure, whereas AI-generated speech tends to have **irregular** phase structure.

B. Complex Vector Spaces \mathbb{C}^N

The set \mathbb{C}^N of all N -tuples of complex numbers forms a vector space over the field \mathbb{C} . A vector $\mathbf{x} \in \mathbb{C}^N$ is written as

$$\mathbf{x} = \begin{pmatrix} x_0 \\ x_1 \\ \vdots \\ x_{N-1} \end{pmatrix}, \quad x_n \in \mathbb{C} \quad (4)$$

Digital audio signals are naturally elements of \mathbb{R}^N in the time domain and \mathbb{C}^N in the frequency domain after FFT transformation. The standard inner product on \mathbb{C}^N is defined as

$$\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{n=0}^{N-1} x_n \overline{y_n} \quad (5)$$

where $\overline{y_n}$ denotes the complex conjugate of y_n . This inner product induces the Euclidean norm.

$$\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} = \sqrt{\sum_{n=0}^{N-1} |x_n|^2} \quad (6)$$

The squared norm $\|\mathbf{x}\|^2$ represents the total energy of the signal.

C. Discrete and Fast Fourier Transform

The Discrete Fourier Transform (DFT) converts a time-domain signal $\mathbf{x} \in \mathbb{C}^N$ to its frequency-domain representation $\mathbf{X} \in \mathbb{C}^N$ via:

$$X_k = \sum_{n=0}^{N-1} x_n e^{-j2\pi kn/N}, \quad k = 0, 1, \dots, N-1 \quad (7)$$

This transformation can be interpreted as computing the inner product of the signal with complex exponential basis vectors at each frequency k . The inverse DFT reconstructs the time-domain signal:

$$x_n = \frac{1}{N} \sum_{k=0}^{N-1} X_k e^{j2\pi kn/N}, \quad n = 0, 1, \dots, N-1 \quad (8)$$

The naive DFT requires $O(N^2)$ complex multiplications. The Fast Fourier Transform (FFT) reduces this to $O(N \log N)$ using a divide-and-conquer strategy. The key insight is that DFT can be decomposed into smaller DFTs of even and odd indexed elements.

Let $\omega_N = e^{-j2\pi/N}$ be the primitive N -th root of unity. For $N = 2^m$, we split the DFT:

$$X_k = \sum_{n=0}^{N-1} x_n \omega_N^{kn} \quad (9)$$

$$= \sum_{n=0}^{N/2-1} x_{2n} \omega_N^{k(2n)} + \sum_{n=0}^{N/2-1} x_{2n+1} \omega_N^{k(2n+1)} \quad (10)$$

$$= \sum_{n=0}^{N/2-1} x_{2n} (\omega_N^2)^{kn} + \omega_N^k \sum_{n=0}^{N/2-1} x_{2n+1} (\omega_N^2)^{kn} \quad (11)$$

Since $\omega_N^2 = \omega_{N/2}$, this becomes

$$X_k = E_k + \omega_N^k O_k \quad (12)$$

where E_k is the DFT of even-indexed elements and O_k is the DFT of odd-indexed elements, both of size $N/2$.

Using the symmetry property $X_{k+N/2} = E_k - \omega_N^k O_k$, we compute both halves with a single recursion:

$$X_k = E_k + \omega_N^k O_k, \quad k = 0, \dots, N/2-1 \quad (13)$$

$$X_{k+N/2} = E_k - \omega_N^k O_k, \quad k = 0, \dots, N/2-1 \quad (14)$$

The recursion bottoms out at $N = 1$, where $X_0 = x_0$. The total complexity satisfies:

$$T(N) = 2T(N/2) + O(N) = O(N \log N) \quad (15)$$

The FFT preserves several important properties crucial for our detection method:

1. Energy Preservation (Parseval's Theorem)

$$\sum_{n=0}^{N-1} |x_n|^2 = \frac{1}{N} \sum_{k=0}^{N-1} |X_k|^2 \quad (16)$$

This ensures that total signal energy is conserved between time and frequency domains.

2. Linearity

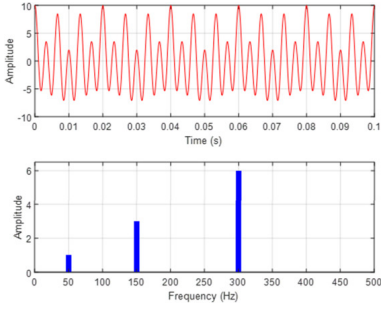


Fig. 2. Time to Frequency Mapping By FFT (Complex Plane)
(Source: <https://www.sciencedirect.com/topics/engineering/fast-fourier-transform>)

$\text{FFT}(\alpha\mathbf{x} + \beta\mathbf{y}) = \alpha\text{FFT}(\mathbf{x}) + \beta\text{FFT}(\mathbf{y})$, allowing us to analyze signal components independently.

3. Symmetry for Real Signals:

When $x_n \in \mathbb{R}$, we have $X_{N-k} = \overline{X_k}$, so only the first $N/2$ coefficients need to be stored.

For our implementation, audio signals are real-valued in the time domain, so we extract magnitude and phase from the positive frequency components $k = 0, \dots, N/2 - 1$.

D. Magnitude and Phase as Complex Algebraic Objects

Each FFT coefficient admits the polar decomposition:

$$X_k = |X_k|e^{j\theta_k} \quad (17)$$

where:

- $|X_k| = \sqrt{\text{Re}(X_k)^2 + \text{Im}(X_k)^2}$ is the spectral magnitude
- $\theta_k = \arctan\left(\frac{\text{Im}(X_k)}{\text{Re}(X_k)}\right)$ is the spectral phase

Phase coherence measures the consistency of phase relationships across frequency bins. For a phase spectrum $\boldsymbol{\theta} = (\theta_0, \theta_1, \dots, \theta_{N-1})$, we define

$$C(\boldsymbol{\theta}) = \left| \frac{1}{N} \sum_{k=0}^{N-1} e^{j\theta_k} \right| \quad (18)$$

This metric ranges from 0 to 1:

- $C = 1$: Perfectly coherent
- $C = 0$: Completely random
- $0 < C < 1$: Partial coherence

The phase velocity (or phase derivative) measures the smoothness of phase progression

$$v_k = \theta_{k+1} - \theta_k \pmod{2\pi} \quad (19)$$

The variance of phase velocity indicates regularity:

$$\sigma_v^2 = \text{Var}(v_k) \quad (20)$$

Low variance suggests smooth and natural phase evolution, while high variance indicates abrupt phase changes typical of synthesis artifacts.

In our implementation, we quantify phase smoothness using the mean absolute phase gradient:

$$\bar{v} = \frac{1}{N-1} \sum_{k=0}^{N-2} |v_k| \quad (21)$$

where smaller values of \bar{v} indicate smoother phase transitions characteristic of natural speech, while larger values suggest the discontinuous phase patterns often present in synthesized audio.

Modern AI voice synthesis (neural vocoders, GANs) primarily optimizes for perceptually accurate magnitude spectra because human hearing is more sensitive to magnitude than phase. However, these models often fail to produce coherent phase structure because phase reconstruction is mathematically ill-conditioned.

E. Spectral Entropy and Energy Distribution

Spectral entropy quantifies the concentration of energy across the frequency spectrum. For a magnitude spectrum $\mathbf{M} = (|X_0|, |X_1|, \dots, |X_{N-1}|)$, we first normalize to obtain a probability distribution

$$p_k = \frac{|X_k|}{\sum_{i=0}^{N-1} |X_i|} \quad (22)$$

The spectral entropy is then defined as

$$H(\mathbf{M}) = - \sum_{k=0}^{N-1} p_k \log p_k \quad (23)$$

This metric characterizes the distribution of spectral energy

- **Low entropy**: Energy concentrated in few frequency bins (typical of human speech)
- **High entropy**: Energy spread uniformly across frequencies (may indicate synthesis artifacts)

The spectral L2 norm

$$\|\mathbf{M}\|_2 = \sqrt{\sum_{k=0}^{N-1} |X_k|^2} \quad (24)$$

represents the total energy in the signal and serves as a normalization factor for geometric distance computations.

F. Window-Based Phase Coherence Analysis

While global phase coherence provides an overall measure, local phase patterns can reveal subtle artifacts. We employ a sliding window approach to compute localized coherence. For a window size w , the local phase coherence at position i is

$$C_i^{(w)} = \frac{1}{w} \left| \sum_{k=i}^{i+w-1} e^{j\theta_k} \right| \quad (25)$$

The overall phase coherence is then the mean of all window coherences:

$$C_{\text{window}} = \frac{1}{N-w} \sum_{i=0}^{N-w-1} C_i^{(w)} \quad (26)$$

This windowed approach offers several advantages:

- Captures local phase consistency patterns
- More robust to isolated phase discontinuities
- Reveals frequency-dependent phase artifacts common in AI synthesis

The window size w is typically chosen based on the expected correlation length of natural speech phase patterns.

G. Data Clustering

In our framework, audio samples form clusters in the complex feature space. Let $\mathcal{H} = \{\mathbf{h}_1, \dots, \mathbf{h}_M\}$ be the set of feature vectors from human speech samples, with centroid

$$\mu_H = \frac{1}{M} \sum_{i=1}^M \mathbf{h}_i \quad (27)$$

Similarly, AI-generated samples form a cluster \mathcal{A} with centroid μ_A . The optimal decision boundary lies at the midpoint between centroids:

$$\tau = \frac{\mu_H + \mu_A}{2} \quad (28)$$

For a test sample with phase coherence C , the classification rule is:

$$\text{decision}(C) = \begin{cases} \text{HUMAN} & \text{if } C > \tau \\ \text{AI-GENERATED} & \text{if } C \leq \tau \end{cases} \quad (29)$$

where τ is the optimal threshold computed from the training data statistics.

H. Mahalanobis-Like Distance and Confidence Estimation

To quantify classification confidence, we compute standardized distances from the test sample to each cluster centroid. Given the phase coherence C of a test sample, and the statistics (μ_H, σ_H) and (μ_A, σ_A) from the human and AI training sets respectively, we define:

$$d_H(C) = \frac{|C - \mu_H|}{\sigma_H + \epsilon} \quad (30)$$

$$d_A(C) = \frac{|C - \mu_A|}{\sigma_A + \epsilon} \quad (31)$$

where ϵ is a small constant to prevent division by zero. These distances are analogous to Mahalanobis distances in one dimension, accounting for the variance of each class.

The confidence score is computed from the relative distances:

$$\text{confidence} = 1 - \frac{\min(d_H, d_A)}{\max(d_H, d_A) + \epsilon} \quad (32)$$

This confidence metric has the following properties:

- confidence $\in [0, 1]$
- High confidence (≈ 1): Test sample is much closer to one cluster than the other
- Low confidence (≈ 0): Test sample is equidistant from both clusters (ambiguous case)

The predicted class is the one with smaller distance:

$$\text{prediction} = \begin{cases} \text{HUMAN} & \text{if } d_H < d_A \\ \text{AI-GENERATED} & \text{otherwise} \end{cases} \quad (33)$$

This geometric approach provides both a classification decision and a measure of certainty, essential for practical deployment where uncertain predictions may require human review.

III. METHODOLOGY

A. System Model

The detection system operates through a four-stage pipeline. Those are audio preprocessing, frequency domain transformation, feature extraction, and threshold-based classification. The core hypothesis is that human speech exhibits distinct phase relationships in the frequency domain that differ systematically from AI-generated audio. These differences arise from fundamental mechanisms that human speech originates from physical vocal cord vibrations and acoustic resonance, while AI-generated speech is synthesized digitally, potentially introducing artifacts in phase structure.

B. Signal Processing

The audio signal is first loaded and normalized to a consistent sampling rate and amplitude range. Stereo recordings are converted to mono by averaging channels. The preprocessed time-domain signal then undergoes Fast Fourier Transform (FFT), converting it from amplitude-over-time representation to amplitude-and-phase-over-frequency representation.

The FFT produces a complex-valued vector where each element corresponds to a frequency component. The magnitude represents energy at that frequency, while the phase angle captures timing relationships. The system retains only positive frequency components, exploiting the symmetry property of real-valued signals. This frequency domain representation becomes the foundation for all subsequent analysis.

C. Feature Extraction

From the FFT output, the system extracts three key features characterizing the signal's geometric properties. First, phase coherence serves as the primary discriminative feature. They are used to quantify how consistently phase transitions across adjacent frequency bins. The computation analyzes phase values in sliding windows across the spectrum. High coherence indicates regular phase transitions typical of human speech. On the other hand, low coherence suggests random or artificial patterns common in nonhuman speech.

Second, phase velocity measures the rate of phase change across frequency. Human speech typically exhibits smoother phase curves due to continuous physical sound production, while AI-generated speech may show abrupt transitions.

Third, spectral entropy quantifies energy distribution across frequencies. Low entropy indicates concentration in specific bands, while high entropy suggests uniform distribution.

D. Classification Strategy

The classification employs a geometric threshold-based approach rather than machine learning. Before classification, the system computes reference statistics from known human and AI-generated speech samples. For each dataset, phase coherence values are extracted from multiple files, and statistical measures (mean and standard deviation) are computed. The decision threshold is established at the midpoint between the human and AI mean phase coherence values.

For a test audio sample, the system extracts its phase coherence and applies a simple decision rule. If phase coherence exceeds the threshold, it is classified as human. Otherwise, it is classified as nonhuman. To quantify confidence, the system computes normalized geometric distances from the test sample to each class reference, similar to Mahalanobis distance. The distance to each class measures how many standard deviations the test sample deviates from that class's mean. Confidence reflects how much closer the sample is to its predicted class compared to the alternative.

This threshold-based strategy offers immediate interpretability with a clear decision boundary, minimal calibration requirements, and avoidance of overfitting risks. The approach prioritizes mathematical transparency and geometric reasoning over complex learned parameters, making the classification process fully explainable through linear algebra principles.

IV. IMPLEMENTATION

V. CASE ANALYSIS

VI. CONCLUSION

REFERENCES

- [1] Ni, Y., et al. (2024). "A Deepfake Detection Algorithm Based on Fourier Transform of Biological Signal." Tech Science Press. (Accessed on December 16, 2025) <https://www.techscience.com/cmc/v79n3/57116>
- [2] Mai, K. T., Bray, S., Davies, T., and Griffin, L. D. Warning: Humans cannot reliably detect speech deepfakes. PLOS ONE, 18(8):1–20, 2023. (Accessed on December 19, 2025) <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0285333>
- [3] Munir, R. (2025). IF2123 Aljabar Linier dan Geometri - Semester I Tahun 2025/2026: Aljabar Kompleks (Update 2024). (Accessed on December 19, 2025) <https://informatika.stei.itb.ac.id/~rinaldi.munir/AljabarGeometri/2025-2026/Algeo-25-Aljabar-Kompleks-2025.pdf>
- [4] Christensen, M. G. (2019). Introduction to Audio Processing. Springer Nature Switzerland AG. (Accessed on December 19, 2025) <https://doi.org/10.1007/978-3-030-11781-8>
- [5] Burrus, C.S. (1989). Algorithms for Discrete Fourier Transform and Convolution. Springer Science+Business Media New York. (Accessed on December 20, 2025). <https://link.springer.com/book/10.1007/978-1-4757-3854-4>
- [6] Uppada, S. K. (2014). Centroid Based Clustering Algorithms- A Clarion Study. International Journal of Computer Science and Information Technologies, Vol. 5 (6), 7309-7313. (Accessed on December 20, 2025). <https://www.ijcsit.com/docs/Volume%205/vol5issue06/ijcsit2014050688.pdf>
- [7] S. Nordholm, S., Togneri, R., Toh, A. M. (2005). Spectral Entropy as Speech Features for Speech Recognition. PEECS. (Accessed on December 20, 2025). https://www.academia.edu/download/95656320/Spectral_entropy_as_speech_features_for_speech_rec.pdf
- [8] Xiang, S., et al. (2008). Learning a Mahalanobis Distance Metric for Data Clustering and Classification. Tsinghua National Laboratory for Information Science and Technology (TNList). (Accessed on December 20, 2025). <https://doi.org/10.1016/j.jesp.2017.09.011>

STATEMENT

Hereby I declare that this paper that I have written is my own work, not a reproduction or translation of someone else's work and not plagiarized.

Bandung, 24 December 2025



Kevin Wirya Valerian
13524019