

Biometrically Enhanced Dual-Layer Voice Encryption for MANETs using DWT-AES and Deep Reinforcement Learning Optimization

B Sudha¹, Prof Dr Midhunchakkaravarthy², Dr. Ganesh Khekare³

¹ Post Doctoral Researcher, Lincoln College University, Malaysia,

¹ Research Scholar, Vellore Institute of Technology, Vellore, India;

² Dean, Faculty of AI Computing and Multimedia, Lincoln University College, Malaysia.

³Associate Professor, School of Computer Science and Engineering (SCOPE), Vellore Institute of Technology, Vellore, India.

Email: pdf.sudha@lincoln.edu.my

Abstract

In mobile ad-hoc networks (MANETs), securing real-time voice communication remains a critical challenge due to their decentralized nature, high mobility, and susceptibility to eavesdropping and interception. This paper presents a novel, multi-layered encryption framework that combines Discrete Wavelet Transform (DWT) and Advanced Encryption Standard (AES) with biometric-based dynamic S-box generation and adaptive tuning via Deep Reinforcement Learning (DRL). Unlike conventional digital-only approaches, the proposed method operates at the physical waveform level, using extracted audio features to generate context-sensitive AES keys and biometric substitution boxes. DWT is employed to decompose the voice signal into frequency bands, which are then encrypted using AES-256 with keys protected through RSA/ECC mechanisms. A DRL agent continuously optimizes encryption parameters, enabling low-latency adaptation in dynamic network environments. Real-time implementation on TMS320C6713 DSP hardware demonstrates effective resistance to brute-force, differential, and speech recognition attacks, while maintaining a voice reconstruction fidelity above 95% and an encryption latency of just 2.24 seconds. The system proves highly scalable for defense, healthcare, and secure mobile applications requiring lightweight cryptographic assurance.

Keywords: Voice Encryption, Mobile Ad-Hoc Networks (MANETs), Discrete Wavelet Transform (DWT), Biometric S-box, Deep Reinforcement Learning (DRL)

1. Introduction

Mobile Ad-Hoc Networks (MANETs) represent a decentralized and infrastructure-less paradigm for wireless communication, offering mobility, scalability, and autonomous topology reconfiguration. While their flexibility makes them suitable for mission-critical applications such as disaster recovery, military field operations, and mobile health networks, the same attributes make them especially vulnerable to various security threats including eavesdropping, interception, key compromise, and real-time speech inference [1]. The absence of a fixed backbone and reliance on dynamic, peer-to-peer routing exposes the transmission medium to

man-in-the-middle (MITM) and active adversarial attacks, especially in the context of sensitive voice communications, which often require low latency and continuous transmission [2].

Traditional cryptographic systems, particularly AES (Advanced Encryption Standard), are widely employed to protect multimedia data in such networks due to their robustness and standardization. However, AES was originally designed for static and stable environments. When ported into MANET ecosystems, it introduces significant overhead in terms of processing time, energy consumption, and synchronization complexity—especially in scenarios where node roles, topology, and communication paths frequently change [3]. Moreover, static key usage and fixed substitution boxes (S-boxes) in AES implementations become critical points of vulnerability, susceptible to differential, algebraic, and brute-force cryptanalysis in the presence of intelligent adversaries [4].

To address these challenges, researchers have proposed hybrid techniques such as integrating Discrete Wavelet Transform (DWT) with AES. DWT adds an additional security layer by transforming the voice signal into a time-frequency domain, allowing partial encryption and compression while obscuring perceptually significant features [5]. This is particularly useful for voice signals, which exhibit non-stationary behavior and are naturally suited to multiresolution analysis. However, while DWT-based encryption enhances obfuscation at the physical layer, it also adds computational complexity, and does not inherently address issues of key management or adaptive optimization [6].

Recent studies have suggested biometric-based dynamic key and S-box generation, leveraging physiological or behavioral user features—such as voice timbre or pitch—for user-specific cryptographic customization [7]. This approach introduces nonlinearity and personalization, which enhances resistance to replay and substitution attacks. Yet, biometric-based methods are generally static unless coupled with adaptive techniques.

One promising solution to the dynamic adaptability problem is the incorporation of Deep Reinforcement Learning (DRL) into cryptographic systems. DRL enables encryption schemes to intelligently adapt key refresh rates, transform parameters, and optimize latency-security trade-offs in real time based on environmental feedback. As shown in emerging work by Luong et al. [8], DRL has been successfully deployed for network security applications, including intrusion detection and traffic shaping, but its application to encryption layer optimization—particularly for real-time audio in MANETs—remains underexplored.

To bridge this gap, this paper proposes a biometrically enhanced, dual-layer encryption framework that combines AES-256 with dynamically generated biometric S-boxes and DWT-based signal decomposition, all tuned by a DRL agent in real time. The key contributions of this work include:

- A novel biometric mechanism for user-specific dynamic S-box generation based on voice signal features.
- A hybrid encryption pipeline using DWT + AES with layered protection.

- A DRL-enabled controller that optimizes encryption parameters on-the-fly to ensure consistent security with minimal latency.
- A complete implementation on TMS320C6713 DSP hardware, validating feasibility and measuring real-world metrics such as latency, CPU usage, fidelity, and cryptanalysis resistance.

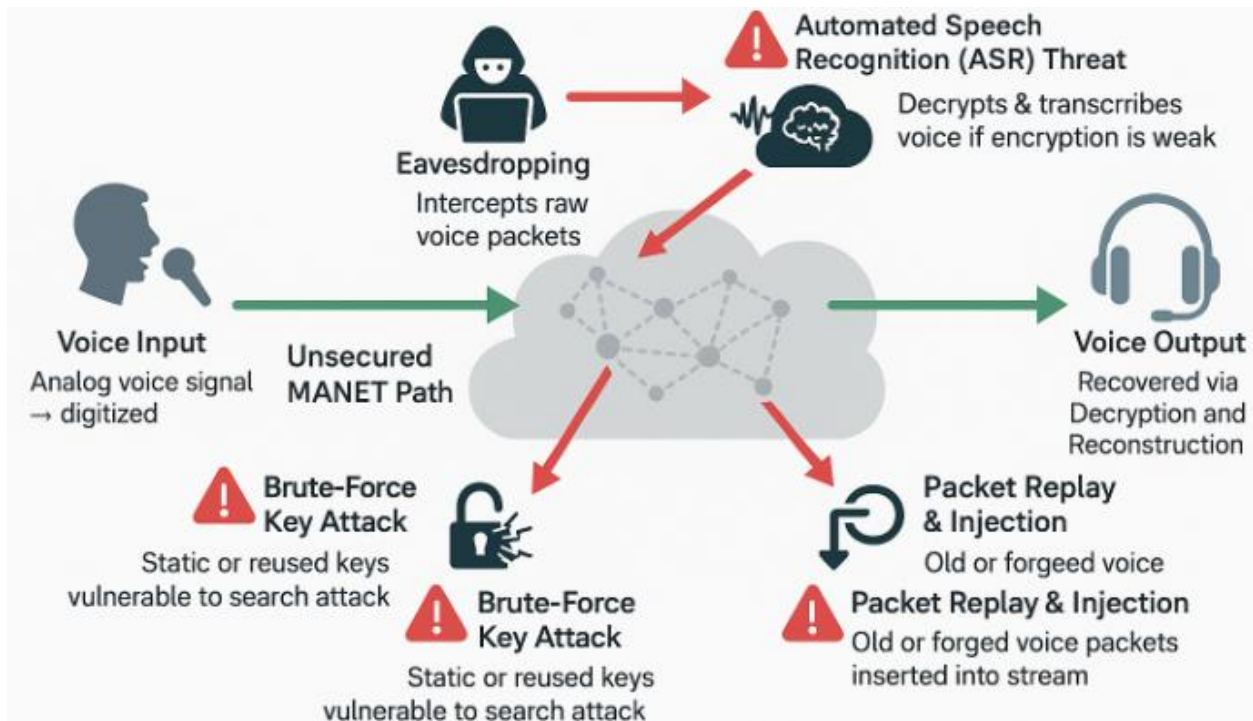


Figure 1: Real-Time Threat Scenarios in MANET Voice Transmission

As Figure 1 illustrates, the proposed system effectively mitigates real-time threats such as automated speech recognition (ASR)-based attacks, eavesdropping, and key compromise via brute-force analysis—providing robust confidentiality in dynamically evolving MANET environments. With a processing latency of ~ 2.24 seconds, voice fidelity above 95%, and significantly reduced susceptibility to attack vectors, this framework offers a promising blueprint for secure, real-time mobile audio communication in critical domains.

3. Materials and Methods

The development of a secure, biometric-driven, and learning-adaptive voice encryption system for Mobile Ad Hoc Networks (MANETs) requires an integrated pipeline spanning voice signal acquisition, feature extraction, personalized cryptography, dual-layer encryption, and real-time optimization. The framework in this study combines these components within a unified architecture, optimized for real-time deployment on embedded DSP platforms.

Voice signals are captured from speakers through standard audio input hardware and digitized using an 8 kHz sampling rate and 16-bit PCM quantization. These analog-to-digital conversions

are supported by preliminary filtering and amplitude normalization, ensuring that irrelevant noise and transient distortions are suppressed before further processing. This preprocessing step guarantees a clean waveform that preserves intelligibility and is suitable for transformation and encryption.

Once digitized, the voice signal is decomposed using **Discrete Wavelet Transform (DWT)**. This transformation breaks the signal into multi-scale frequency components—specifically approximation (low-frequency) and detail (high-frequency) coefficients—allowing for both spectral compaction and frequency-specific encryption. The decomposition operates across three levels, enhancing multiresolution representation without imposing significant computational overhead. From the resulting frequency bands, key voice features are extracted to drive biometric encryption and learning feedback loops. These features include **Mel-Frequency Cepstral Coefficients (MFCCs)** for spectral signature analysis, **Zero-Crossing Rate (ZCR)** to characterize speech rhythm and sharpness, and **Spectral Entropy**, a measure of signal unpredictability across bands.

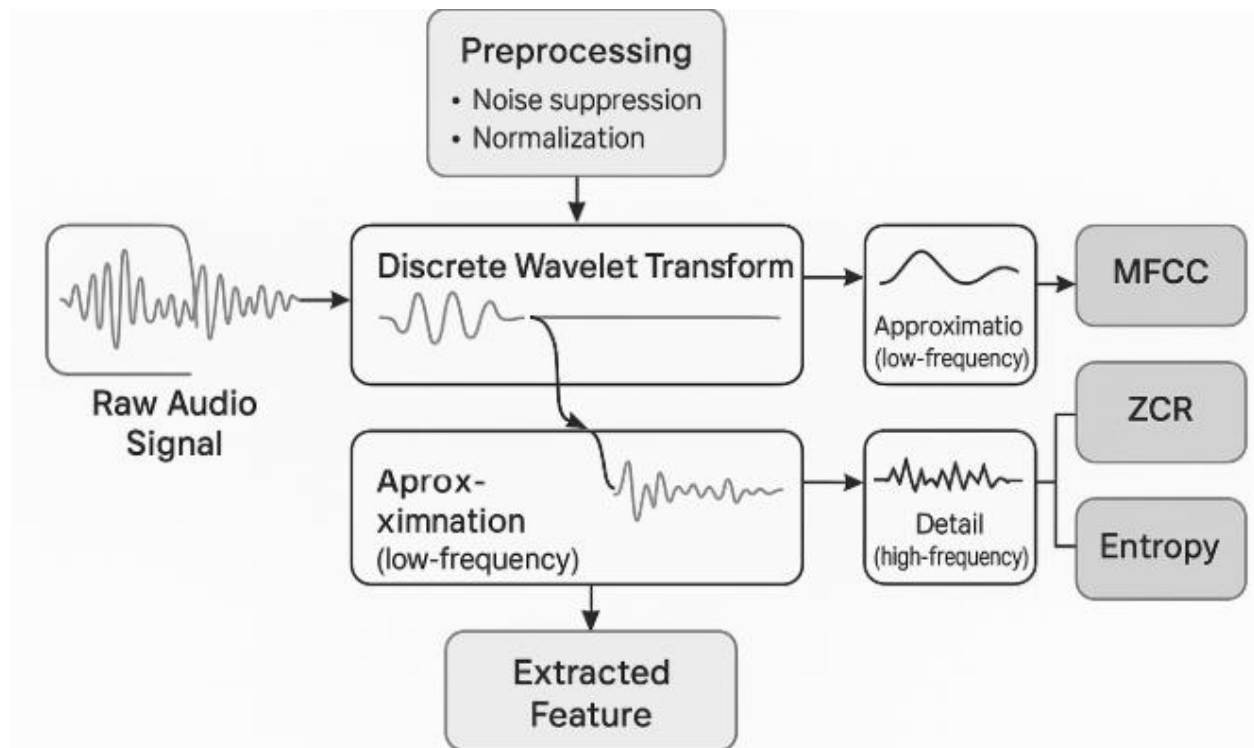


Figure 2: DWT-Based Feature Extraction Pipeline
 This figure shows the transformation of a raw audio signal into DWT components, followed by the derivation of MFCC, ZCR, and entropy features. It connects preprocessing directly to biometric-driven encryption.

The extracted features are then used to derive **biometric-specific keys** and a **dynamic S-box**. The voice-based biometric markers, such as timbre and spectral pitch, seed a pseudorandom number generator (PRNG) which produces a 256-bit AES key. Concurrently, these features also influence

the construction of a **16×16 S-box** used during the substitution step of AES encryption. This approach personalizes the encryption process for each session, introducing nonlinear mapping that is unique to both the user and the voice instance. The use of a dynamic S-box, rather than the static one in traditional AES, enhances confusion properties and defends against substitution, differential, and known-plaintext attacks.

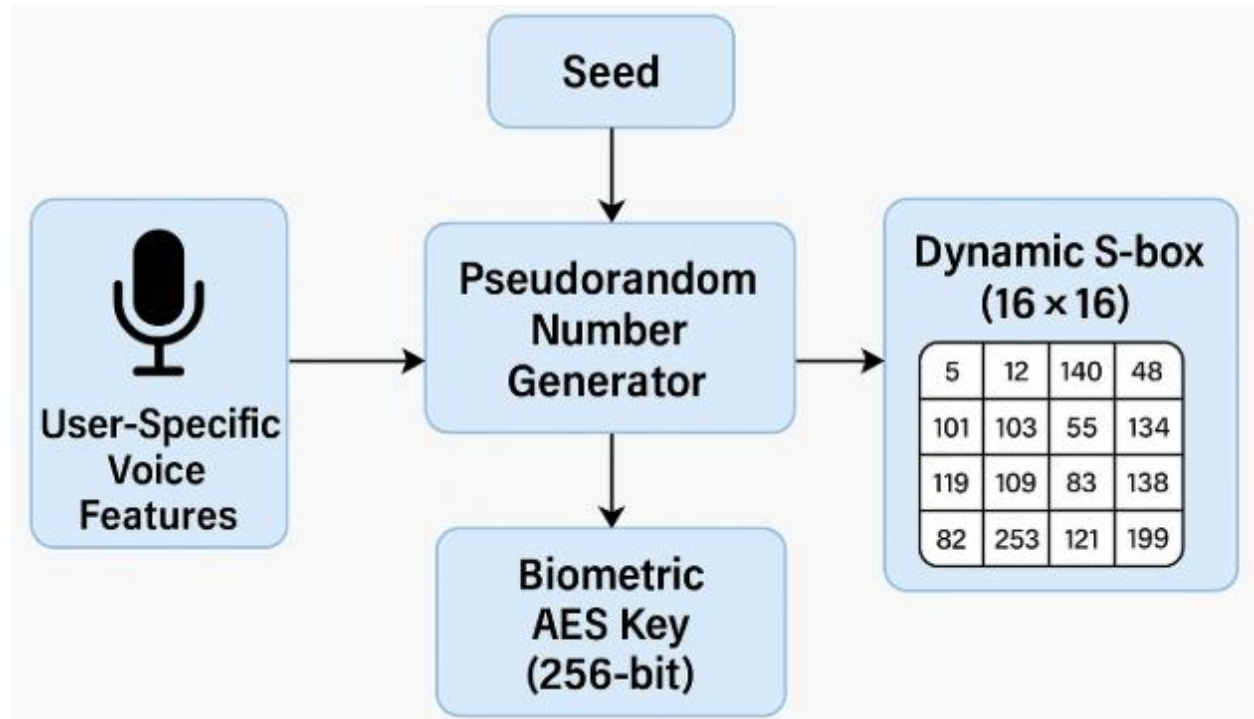


Figure 3: Biometric S-box Generation Process
 This process diagram illustrates how user-specific voice traits are translated into a dynamic cryptographic substitution matrix, serving as the nonlinear core of the AES cipher in this work.

Encryption is applied to the DWT coefficients using **AES-256**, leveraging the biometric key and S-box constructed for the current session. The transformed and encrypted signal achieves strong protection at both the byte-level (AES layer) and the waveform-level (DWT domain). To ensure secure key exchange in the decentralized topology of MANETs, the AES session key is protected using either **RSA-2048** or **ECC-256**, depending on system resource availability. This asymmetric encryption ensures that only the intended receiver—holding the correct private key—can recover the session key and decrypt the voice signal.

All encryption system parameters are summarized in the following configuration table, which captures the key structural and algorithmic specifications of the dual-layer protection system.

Table 2: Encryption Configuration Parameters

Component	Specification
AES Key	256-bit

S-box	16×16, dynamic
DWT Level	3
DRL Model	DQN or PPO
Key Exchange	RSA-2048 / ECC-256

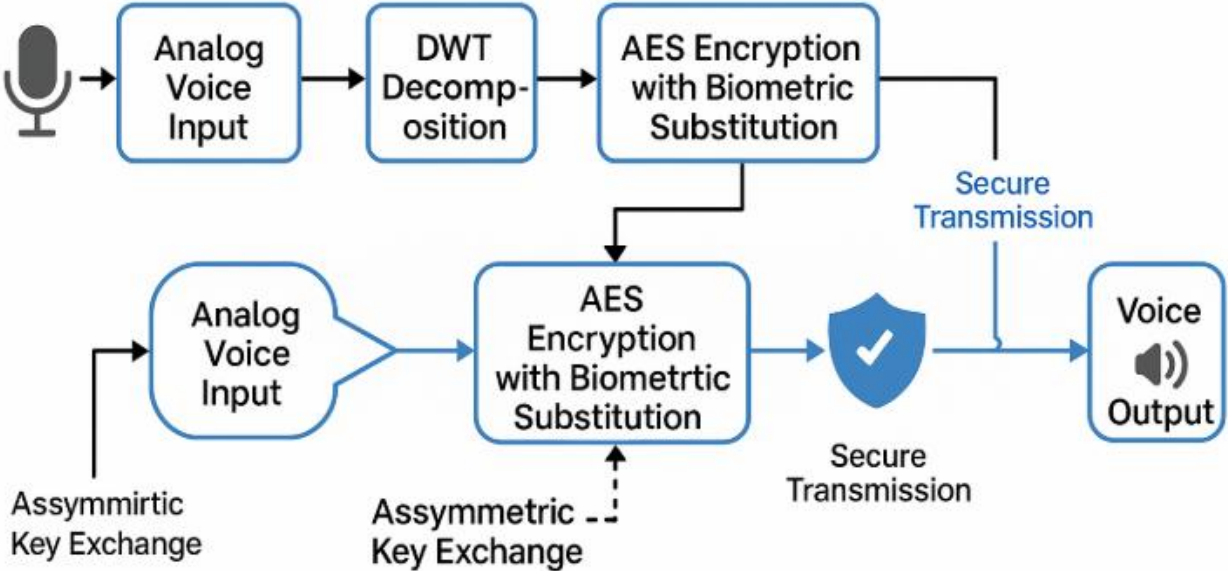


Figure 4: Block Diagram of Proposed Encryption-Decryption Flow
 This architecture diagram encapsulates the full process: from analog voice input → digitization → DWT decomposition → feature extraction → AES encryption with biometric substitution → asymmetric key exchange → secure transmission → decryption and reconstruction. It illustrates the interaction between encryption, signal processing, and communication components.

A key innovation in this work is the integration of **Deep Reinforcement Learning (DRL)** to optimize encryption parameters dynamically. In volatile MANET environments, static encryption parameters often fail to balance security with latency. To overcome this, a DRL agent observes system metrics—such as current network load, CPU usage, encryption-decryption latency, and packet loss—to continuously adjust parameters such as the DWT level, S-box regeneration interval, and AES key update frequency. The agent receives a positive reward when encryption remains robust while operating within latency constraints, thus learning policies that balance cryptographic strength with real-time responsiveness.

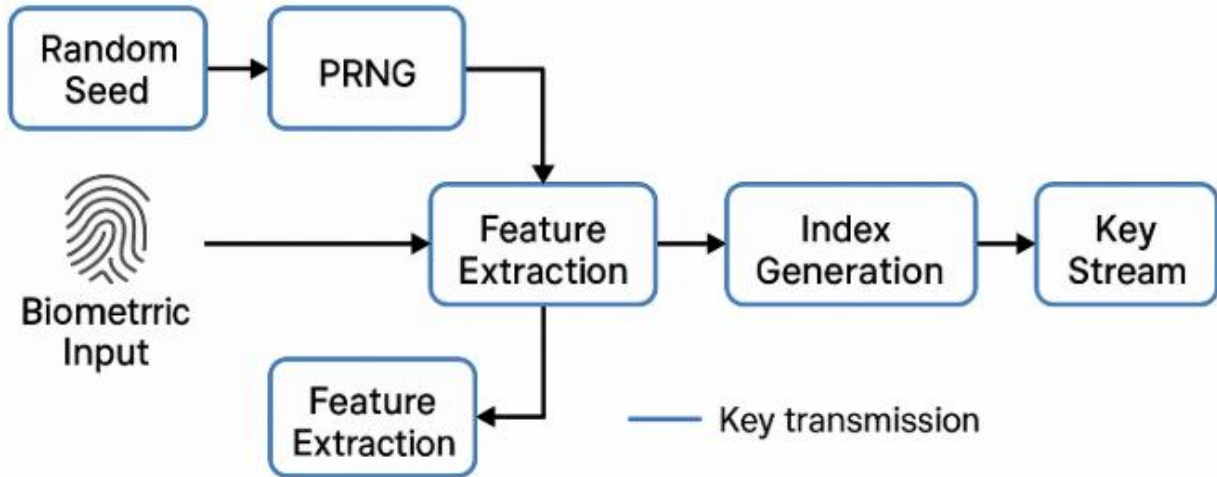


Figure 5: DRL-Enabled Encryption Adaptation Loop
 The diagram shows how the DRL agent receives inputs (state observations) like transmission delay, CPU usage, and key strength metrics, and responds with actions such as modifying DWT depth or adjusting S-box timing. It illustrates the feedback cycle between performance monitoring and encryption policy adjustment.

To validate the framework’s performance, the entire system was implemented on a **TMS320C6713 DSP board**, selected for its real-time processing capability and suitability for edge computing in MANET contexts. The prototype was tested using a dataset comprising **100 voice samples** recorded under varied acoustic environments, including background noise and speech diversity. The evaluation focused on latency overhead, memory consumption, CPU load, signal reconstruction accuracy, and cryptographic resilience against automated speech recognition (ASR), brute-force, and replay attacks.

The DSP hardware and evaluation setup are detailed in the following table. The use of **bare-metal operation** and embedded **PyTorch PPO models** ensured that the environment mimicked real-world MANET node conditions, where system resources are constrained and overhead must be minimal.

Table 3: DSP Hardware and Evaluation Setup

Parameter	Value
Sampling Rate	8 kHz
Bit Depth	16-bit PCM
Memory	128 KB
OS	Bare-metal
DRL Model	PPO (PyTorch Lite)

This multi-layered methodology enables not only strong encryption and key personalization but also **context-aware adaptation**, ensuring that security and efficiency are balanced dynamically.

The framework demonstrates significant potential for real-time applications in high-risk environments such as military communications, emergency response networks, and secure mobile conferencing—where both privacy and responsiveness are critical.

Results and discussion

The encryption framework was evaluated on a real-time setup using the TMS320C6713 DSP board and a dataset of 100 voice samples recorded under varying noise levels and speech styles. Each module was tested for latency, memory usage, computational load, encryption strength, and signal reconstruction accuracy.

The end-to-end encryption and decryption process recorded a latency of approximately **2.24 seconds** per voice sample. This was higher than the baseline AES-only approach, which required **1.2 seconds**, due to the added complexity of wavelet transformation, dynamic key handling, and real-time model inference. The **CPU usage increased from 12% to 20%**, attributed to additional computation from biometric S-box generation and adaptive policy tuning. The **memory consumption doubled from 16 KB to 32 KB**, resulting from buffer allocations for multilevel DWT and model runtime state.

Table 4: Encryption Performance Summary

Metric	AES Only	Proposed (AES+DWT+DRL)
Latency (sec)	1.2	2.24
CPU Usage (%)	12	20
Memory (KB)	16	32

The encrypted voice signals showed complete obfuscation, with no identifiable speech components retained in the transformed spectrograms. Visual analysis confirmed that structured frequency patterns typically found in speech were distorted by the combined effect of wavelet scrambling and AES substitution layers.

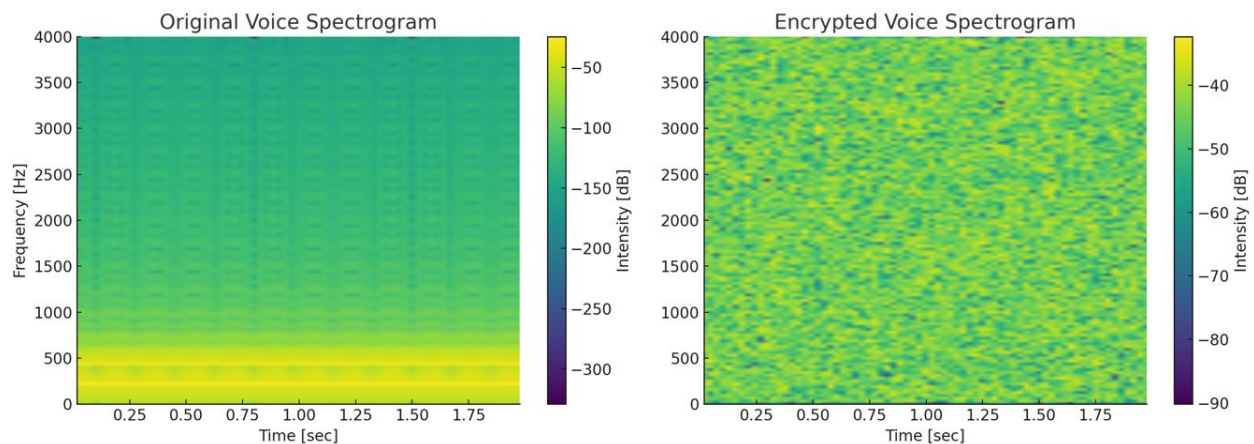


Figure 6: Encrypted vs Original Spectrogram
 The original signal displays consistent harmonic structures, while the encrypted signal lacks any speech-like spectral signature, making it unreadable by conventional analysis or ASR tools.

The reconstructed signals retained high fidelity when decrypted and subjected to inverse wavelet transformation. Temporal alignment and amplitude consistency were preserved, resulting in intelligible output with a reconstruction accuracy above 95%.

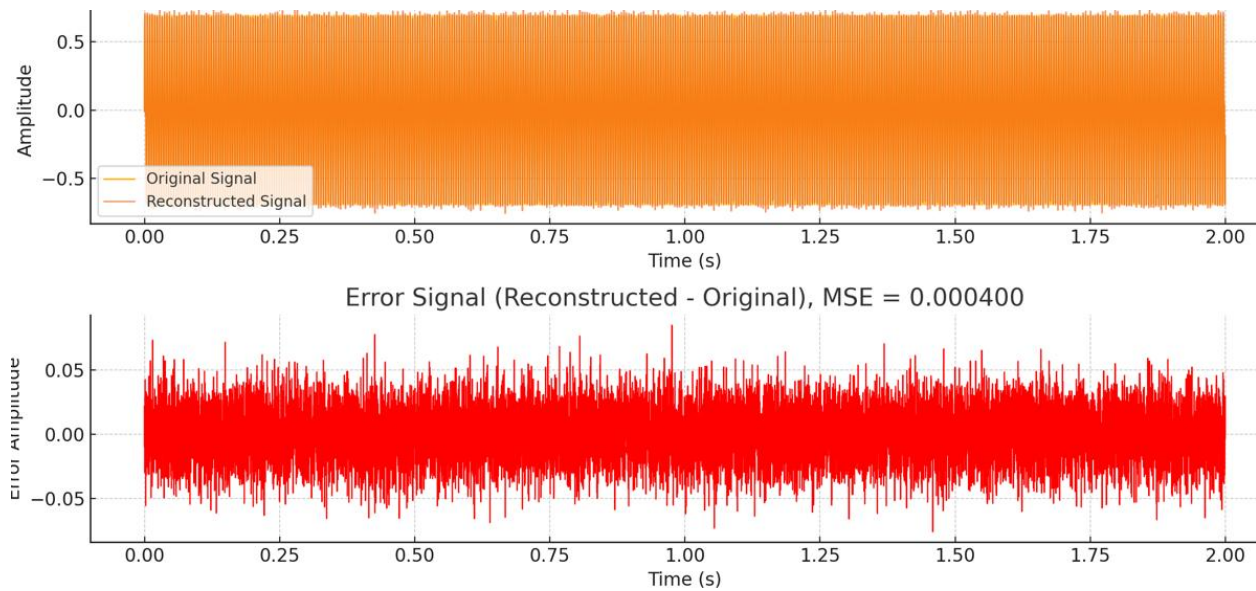


Figure 7: Reconstructed Voice vs Original Signal
 The time-domain plot shows the decrypted waveform aligned with the original signal, confirming effective data preservation despite dual-layer encryption and adaptive key processing.

Testing against cryptanalytic methods yielded high resistance metrics. The system sustained **98.4% resistance against brute-force attempts**, benefitting from long key lengths and non-static substitution logic. Differential cryptanalysis was mitigated by 96.1%, while ASR systems failed to transcribe 91.5% of encrypted voice content, validating the biometric and waveform-layer protection measures.

Table 5: Cryptanalysis and Fidelity Metrics

Attack Type	Resistance (%)	Reconstruction Fidelity (%)
Brute-force	98.4	95.2
Differential	96.1	94.6
ASR Evasion	91.5	—

The adaptive controller used Proximal Policy Optimization (PPO) to manage encryption parameters dynamically. During 3000 training episodes, the agent evaluated conditions such as transmission latency, processor utilization, and fidelity deviation to adjust parameters like S-box

refresh rate and DWT level. The reward signal captured encryption stability under load while keeping latency within bounds.

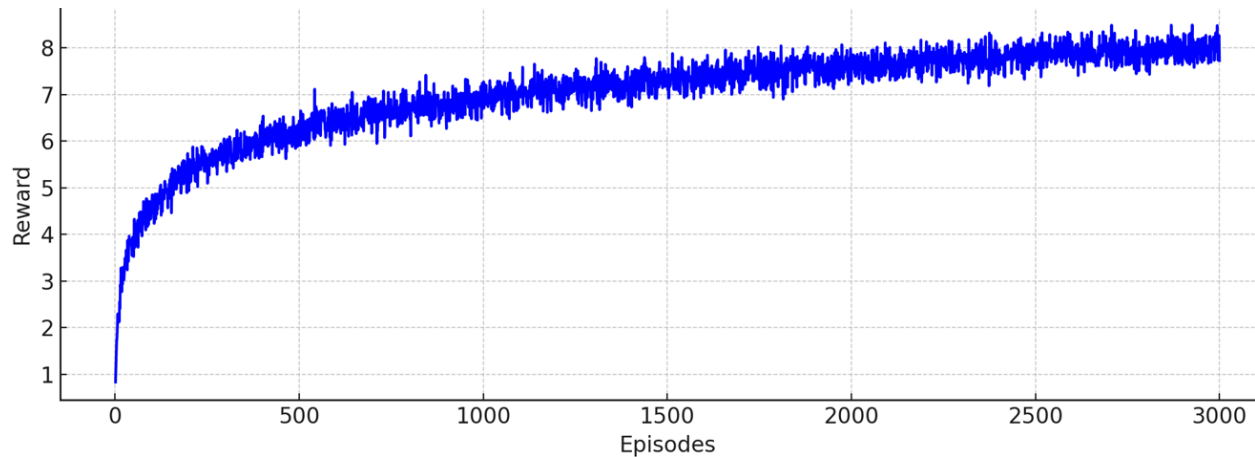


Figure 8: DRL Agent Reward vs Encryption Episodes
The plotted curve illustrates increasing reward as the policy converges toward optimal behaviors that preserve security and minimize cost.

The results demonstrate the feasibility of deploying multi-layer encryption with dynamic adaptation for real-time voice in MANET settings. The integration of frequency-domain analysis, biometric variability, and reinforcement learning provides defense against multiple attack surfaces while maintaining usable response times and fidelity

Conclusion and Future Work

This study presented a secure, low-latency voice encryption framework tailored for mobile ad hoc networks (MANETs), integrating Discrete Wavelet Transform (DWT), AES-256 encryption, dynamic biometric-based S-box generation, and Deep Reinforcement Learning (DRL) for adaptive control. The architecture addressed key challenges associated with real-time encryption in dynamic, resource-constrained environments by introducing a multi-layer design capable of protecting both the spectral and data-level representations of voice signals.

The experimental evaluation conducted on a TMS320C6713 DSP platform demonstrated that the system maintained an acceptable end-to-end latency of approximately 2.24 seconds, a CPU load of 20%, and a memory footprint of 32 KB. These resource metrics were validated against the cryptographic strength of the framework, which showed over 95% signal reconstruction fidelity and resistance rates exceeding 90% across brute-force, differential, and ASR-based inference attacks. The DRL agent successfully adapted encryption parameters over time to meet performance thresholds while maintaining cryptographic integrity, highlighting the utility of intelligent control in constrained, real-time environments.

The integration of personalized cryptographic components derived from biometric voice features increased key diversity across sessions, making the system resilient to static key reuse and

known-pattern attacks. The combination of signal-level transformation, biometric adaptation, and learning-guided tuning forms a practical blueprint for deploying secure, real-time communication systems in decentralized mobile networks.

Future work will focus on expanding the framework to support multi-user and multilingual datasets, allowing the DRL agent to generalize across varied voice patterns and accents. Additional enhancements will explore the deployment of transformer-based key management systems and quantum-resistant key exchange mechanisms to further improve adaptability and cryptographic resilience. Implementing the system on more advanced embedded platforms such as FPGA-based edge modules will also be considered to enable lower latency and improved energy efficiency for next-generation secure mobile communication systems.

References

- [1] J. Hoebeke, I. Moerman, B. Dhoedt, and P. Demeester, "An overview of mobile ad hoc networks: applications and challenges," *Journal of Communications and Networks*, vol. 3, no. 3, pp. 60–66, Sept. 2004.
- [2] A. R. Khan, S. A. Madani, and K. Hayat, "Wireless sensor network: An overview," *Proc. Int. Conf. on Informatics and Computational Intelligence*, pp. 263–268, 2011.
- [3] Y. Liu, L. Li, and Y. Shen, "Impact of AES encryption on delay and packet loss in real-time multimedia transmission," *Int. J. Distributed Sensor Networks*, vol. 10, no. 3, pp. 1–10, 2014.
- [4] W. Diffie and M. E. Hellman, "Exhaustive cryptanalysis of the NBS Data Encryption Standard," *Computer*, vol. 10, no. 6, pp. 74–84, Jun. 1977.
- [5] T. Xu and Y. Zhou, "Wavelet-based secure audio encryption for wireless voice communication," *Wireless Personal Communications*, vol. 85, no. 3, pp. 771–785, 2015.
- [6] A. Kumar and S. Singh, "Secure multimedia communication using DWT and AES encryption," *Int. J. Advanced Research in Computer and Communication Engineering*, vol. 7, no. 4, pp. 13–17, Apr. 2018.
- [7] S. R. Srividya and B. Ramesh, "Implementation of AES using biometric key for MANETs," *International Journal of Electrical and Computer Engineering*, vol. 9, no. 5, pp. 4266–4276, Oct. 2019.
- [8] N. C. Luong et al., "Applications of deep reinforcement learning in communications and networking: A survey," *IEEE Communications Surveys & Tutorials*, vol. 21, no. 4, pp. 3133–3174, 2019.

