# A THEORETICAL ARCHITECTURE OF VOICEPRINT RECOGNITION FOR NETWORK SECURITY SITUATIONAL AWARENESS

Ping Xia

School of Engineering, Guangzhou College of Technology and Business, Foshan 528138, Guangdong, China. Corresponding Email: 710398795@qq.com

Abstract: This paper proposes a theoretical framework for DenseNet-based voiceprint recognition, which incorporates spectrogram enhancement and adaptive histogram equalization to overcome the limitations of conventional methods in feature extraction robustness under noisy conditions. The framework synergistically combines spectral feature enhancement with DenseNet's dense connectivity, achieving both improved feature discriminability and deep feature reuse through: optimized time-frequency representation via enhanced spectrograms, hierarchical feature propagation enabled by dense blocks. Theoretical analysis confirms the framework's capability to maintain recognition stability against acoustic interference, establishing a novel biometric authentication paradigm for cybersecurity situational awareness systems.

Keywords: Voiceprint recognition; Spectrogram feature enhancement; Histogram equalization; Cybersecurity; Situational awareness

# **1 INTRODUCTION**

With the continuous advancement and increasing sophistication of cyberattack techniques, traditional password-based authentication systems are facing severe challenges. As a non-intrusive biometric identification technology, voiceprint recognition offers multiple advantages including low-cost voice data acquisition, mature technology, low computational complexity of processing algorithms, and the capability for remote authentication, making it an ideal implementation technology for network identity recognition applications. However, the robustness of existing voiceprint recognition systems in complex network environments remains to be addressed. This paper combines spectrogram enhancement with histogram equalization to establish a fusion model integrating voiceprint recognition and cybersecurity situational awareness, aiming to resolve the issues of insufficient feature extraction and inadequate robustness in traditional voiceprint recognition technologies operating in complex network environments, while providing an interpretable theoretical framework for identity authentication scenarios.

# **2 RELATED WORK**

The evolution of voiceprint recognition technology has progressed from traditional statistical models to contemporary deep learning approaches. Early methodologies predominantly employed machine learning algorithms such as Hidden Markov Models (HMM), Gaussian Mixture Models (GMM), GMM-Universal Background Models (GMM-UBM), GMM-Support Vector Machines (GMM-SVM), and i-vector systems [1-3]. While these approaches established foundational frameworks, they exhibit inherent limitations including oversimplified feature representations, limited recognition accuracy, susceptibility to channel variability, and inadequate noise robustness—constraints that significantly impede their efficacy in complex operational environments [4], particularly within cybersecurity situational awareness applications requiring high reliability.

The advent of artificial intelligence has catalyzed a paradigm shift, with deep learning emerging as the predominant research focus for voiceprint recognition [5-7]. In the context of cybersecurity authentication, voice biometrics now offer enhanced verification solutions for situational awareness systems. Liu et al. developed an LSTM-based architecture utilizing spectrogram representations of voiceprints [8], achieving superior text-independent recognition accuracy by capitalizing on LSTM's sequential modeling capabilities. Wang et al. advanced this domain through an end -to-end bidirectional LSTM framework that exploits temporal dependencies in speech sequences [9], demonstrating scalability for large-scale user authentication—a critical requirement for modern cybersecurity infrastructures.

For real-time network security monitoring, Zhao et al. introduced an end-to-end CNN architecture incorporating MFCC feature extraction and Universal Background Modeling [10], effectively mitigating environmental and individual variability. Yan et al. further optimized computational efficiency through a hybrid CNN-LSTM model processing fixed-length spectrograms [11], achieving high accuracy with reduced training iterations. Recent breakthroughs involve ResNet architectures that extract spatiotemporal voiceprint features [12-13], addressing historical challenges in recognition complexity and accuracy while enhancing practical deployment viability. Among existing research achievements, current studies on voiceprint recognition primarily focus on traditional voiceprint feature extraction methods, the implementation of voiceprint recognition using various deep learning approaches, and the application of multimodal fusion authentication technologies. The advancement of these technologies continues to enhance the application value of voiceprint recognition in cybersecurity situational awareness.

# **3** DESIGN OF VOICEPRINT RECOGNITION MODEL BASED ON DENSENET DEEP LEARNING

This paper proposes an end-to-end voiceprint recognition model based on the DenseNet architecture, which achieves efficient mapping from raw speech signals to speaker identity through a hierarchical feature learning mechanism. As illustrated in Figure 1, the system adopts a streaming processing framework that fully leverages DenseNet's advantages in feature reuse and gradient optimization. The integrated enhancement module at the front-end significantly improves the feature representation capability of traditional spectrograms in complex acoustic environments, thereby delivering a more robust identity authentication solution for cybersecurity situational awareness systems.





# **3.1 Spectrogram Generation**

The spectrogram is a graphical representation of speech signals that transforms one-dimensional time-domain data into a three-dimensional image format, dynamically displaying temporal characteristics through the interplay of timevarying frequency components and energy intensity. Color gradients on the spectrogram form distinct texture patterns, which encode substantial speaker-specific biometric features, making this representation particularly suitable for training voiceprint recognition models using deep learning methodologies. As illustrated in Figure 2, the spectrogram generation process consists of three primary stages:

1. Preprocessing: The raw speech signal undergoes pre-emphasis to amplify high-frequency components, compensating for excessive attenuation during signal transmission. The processed signal is then segmented into fixed-duration frames, with each frame subjected to windowing and zero-padding operations to ensure continuity.

2. Power Spectrum Calculation: Each frame is processed through Fast Fourier Transform (FFT) to obtain its frequency spectrum. The magnitude spectrum is squared to derive the power spectrum, followed by logarithmic scaling (log-power) to enhance the dynamic range of spectral features.

3. Spectrogram Construction: The log-power spectra of individual frames are mapped onto a time-frequency coordinate system. Sequential frames are concatenated along the temporal axis to form the final spectrogram - a time-frequency-energy representation that completes the transformation from acoustic waveforms to visual discriminative features. This conversion from time-domain signals to frequency-domain visualizations provides a foundational analytical framework for subsequent deep feature extraction in voiceprint recognition systems.



Figure 2 Conversion Process from Speech Signal to Spectrogram

#### 3.2 Spectrogram Image Enhancement Algorithm

The integration of spectrogram image enhancement algorithms into voiceprint recognition systems facilitates the extraction of salient frequency-domain features from speech signals. By applying histogram equalization techniques,

this approach effectively disperses concentrated noise distributions while mitigating luminance variations caused by interspeaker differences or recording condition disparities, thereby significantly improving spectrogram quality.

Histogram equalization represents a computationally efficient nonlinear transformation method for spectrogram enhancement, operating through grayscale value redistribution to amplify contrast in images with constrained dynamic range. For grayscale spectrogram representations, the histogram provides a quantitative depiction of intensity level distributions, where visual quality exhibits direct correlation with the statistical moments (mean and variance) of grayscale distributions.

In terms of voiceprint feature processing, an improved histogram equalization algorithm is proposed, which has three key theoretical innovations. First, the quantization grading strategy, which establishes a more stable energy adjustment mechanism by converting the traditional continuous grayscale mapping into discrete grade adjustment, thus effectively avoiding the over-enhancement problem in low-energy frequency bands; second, the frequency domain perception mechanism: based on the physical characteristics of speech signals, an adaptive enhancement scheme is designed in the key frequency bands such as resonance peaks, which realizes the differentiated processing of different frequency bands; and third, the dynamic equilibrium design : The optimized balance model of global enhancement and local feature retention is constructed by introducing intelligent adjustment parameters, which theoretically ensures the quality of feature extraction. The algorithm theoretically realizes the balance between computational complexity and feature enhancement effect, and provides a more reliable input basis for subsequent deep feature extraction. From the theoretical analysis, this improved method is particularly suitable for complex acoustic environments where background noise or channel distortion exists.

#### 3.3 DenseNet-based Voiceprint Recognition Network Model

In deep learning networks, the problem of gradient vanishing becomes more and more obvious as the depth of the network increases. Compared to ResNet, DenseNet's algorithm and network structure, although different, are to connect all layers directly to each other under the premise of ensuring maximum information transfer between layers in the network. The difference is that ResNet combines the layers by accumulating the features before passing them to the next layer, while DenseNet combines them by feature connection, establishes the shortest dense connection between all the layers in front and the layers behind, improves the flow of information and gradient between different layers, and makes the network model easy to train. In the DenseNet network structure, in each layer, all the feature maps of the previous prediction layers are used as inputs to the current layer, and the output feature maps are used as inputs to all the later layers, realizing feature reuse. Since the feature map of each layer of DenseNet can be directly used by all subsequent layers, it realizes the reuse of voiceprint features in the whole network model, effectively reduces the number of parameters, and makes the structure of voiceprint recognition network model more concise.

This paper presents a deep neural network model based on an improved densely connected architecture for speech spectrogram feature extraction. The network structure employs a multi-level feature fusion mechanism that establishes cross-layer feature sharing channels, enabling deep integration and efficient utilization of speech features.

In the network initialization phase, large-scale convolutional kernels (7×7) combined with max-pooling operations (3×3) are employed for preliminary feature extraction, with dynamic normalization processing introduced after the convolutional layers to significantly enhance feature representation stability. During the deep feature learning stage, a feature reuse module is designed to directly connect each convolutional layer's output features to all subsequent layers through dense connectivity. This design not only preserves the integrity of low-level features but also achieves cross-layer feature transmission. Specifically, each feature transformation unit adopts a three-stage processing flow of "normalization-activation-convolution," performing local feature extraction through  $3\times3$  convolutional kernels followed by channel-wise feature concatenation. For network optimization, an adaptive feature compression mechanism is introduced, using  $1\times1$  convolutional kernels for dynamic feature dimension adjustment combined with  $2\times2$  average pooling operations for feature map resolution optimization. This design maintains feature representation capability while effectively controlling computational complexity, achieving an optimal balance between feature extraction efficiency and representation capability to provide a reliable deep feature representation solution for speaker recognition tasks. Furthermore, the network architecture incorporates a dual attention mechanism in both temporal and frequency domains during feature fusion, enabling the network to adaptively focus on key feature regions in speech signals. This design significantly improves the network's robustness in complex acoustic environments.

# 4 APPLICATION FRAMEWORK DESIGN OF VOICEPRINT RECOGNITION IN CYBERSECURITY SITUATIONAL AWARENESS

To address the requirements of cybersecurity situational awareness, this study constructs a five-layer architecture system based on voiceprint recognition, achieving closed-loop management from data collection to security response. The framework adopts a modular design, and its system architecture is shown in Figure 3.



Figure 3 Application Framework of Voiceprint Recognition in Cybersecurity Situational Awareness

#### 4.1 Data Acquisition and Processing

The data acquisition layer serves as the physical sensing terminal of the system, adopting a distributed architecture design to realize real-time voice data capture and preprocessing through multi-node collaboration. This layer primarily accomplishes the collection of multi-source voice signals and employs adaptive noise suppression algorithms to eliminate environmental interference, ensuring input quality for subsequent voiceprint feature processing. It achieves digital conversion and standardized processing of voice signals, compensating for high-frequency component attenuation through pre-emphasis filters. This preprocessing pipeline provides high-quality input data for subsequent processing stages.

#### 4.2 Voiceprint Feature Extraction and Enhancement

The feature processing layer employs hybrid signal processing technology, primarily consisting of three processing stages. In the preprocessing stage, frame splitting and windowing operations are performed using Hamming windows to reduce spectral leakage. The feature extraction stage acquires MFCC features through Mel filter banks and cepstral analysis. The feature optimization stage applies an improved histogram equalization algorithm to enhance spectrogram characteristics. This processing flow effectively improves the accuracy of subsequent recognition.

# 4.3 Multimodal Deep Learning Model

The model layer adopts a DenseNet-based multimodal fusion architecture that dynamically integrates voiceprint features, spectrogram features, and behavioral features through a cross-modal attention mechanism. The model employs a spatiotemporal alignment strategy to ensure feature consistency and incorporates an attention weighting mechanism to highlight key features. Its multitask output layer simultaneously accomplishes voiceprint recognition, anomaly detection, and behavior analysis. Theoretical analysis demonstrates three core advantages of this design: dense connections ensure efficient feature reuse and gradient propagation, multimodal fusion enhances feature discriminability, and optimized transition layers effectively control computational complexity while maintaining performance. This architecture is particularly suitable for identity authentication scenarios in cybersecurity situational awareness that demand high real-time performance and robustness.

# 4.4 Situational Awareness Analysis

The situational awareness analysis layer adopts a multi-dimensional security analysis architecture to build a comprehensive threat assessment model by fusing voiceprint biometrics, user behavioral patterns and contextual environment data. Using dynamic risk quantification algorithms, combined with real-time behavioral analysis and historical baseline comparison, it realizes intelligent scoring and warning of threat levels. At the visualization level, it intuitively presents changes in security posture. Through the introduction of adaptive learning mechanism, the system can continuously optimize the detection threshold, effectively identify voice forgery, abnormal access and other security threats, and provide intelligent decision support for network security protection.

#### 4.5 Response and Feedback Mechanisms

The security response layer first performs a risk assessment of detected security events, and based on the assessed threat level (low/medium/high risk), the system triggers step-by-step upgraded security measures. Low-risk threats may be ordinary abnormal behavior, triggering logging and using email notifications for alerts. Medium-risk threats are suspicious activities, such as multiple abnormal logins, potential brute-force break-ins, etc., and will initiate secondary authentication, requiring the user to re-verify their identity through multi-factor authentication. When the secondary authentication fails (Auth Failed), the system automatically escalates the event to High Risk and performs the highest level of response. High Risk threats such as clear attacks, such as malicious code injection and privilege elevation attempts, require immediate isolation of the source of the attack, traceability and forensics (recording voiceprint fingerprints), termination of the current session and IP blackout. the system generates network simulation of the attack samples through confrontation to continuously optimize the robustness of the model, forming a closed-loop updating mechanism. The system adopts a console visualization monitoring interface to support security administrators to grasp real-time changes in the situation and implement manual intervention, forming a human-computer cooperative intelligent defense system. Each response link realizes asynchronous communication through the message bus to ensure the high availability and scalability of the system.



Figure 4 Threat Response Workflow Based on Risk-Based Scoring (RBS)

# **5** CONCLUSION

This study proposes a voiceprint recognition architecture that integrates DenseNet's feature reuse mechanism with spectrogram enhancement technology, which significantly improves the system's feature discrimination capability in complex acoustic environments and provides a reliable technical foundation for dynamic identity authentication in cybersecurity situational awareness. The end-to-end architecture achieves a complete closed-loop process from voice acquisition to threat assessment, demonstrating the application potential of deep feature learning in the field of network security.

# **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

# FUNDING

The project was supported by the following fundings:

1. 2021 Guangdong General Colleges and Universities Young Innovative Talents Project (Project No. 2021KQNCX126);

2. China Association of Higher Education (CAHE) 2024 Higher Education Research Program "Research on the Application of Big Data Analytics and Artificial Intelligence Technology in Cybersecurity Situational Awareness and Assessment for Higher Education Institutions" (Project No. 24XH0205);

3. 2022 Ministry of Education Industry-University Co-operation Collaborative Education Project (Project No. 220605211102737).

# REFERENCES

- [1] Alam M J, Kenny P, Ouellet P, et al. Multi-task learning for speaker verification and antispoofing using Gaussian mixture models. IEEE/ACM Transactions on Audio, Speech, and Language Processing, 2020, 28: 1696-1709.
- [2] Villalba J, Chen N, Snyder D, et al. State-of-the-art speaker recognition with neural network embeddings in NIST SRE18 and Speakers in the Wild evaluations. Computer Speech & Language, 2021, 60: 101026.
- [3] Ferrer L, McLaren M, Lawson A. Probabilistic linear discriminant analysis with vector embeddings for speaker verification. IEEE Journal of Selected Topics in Signal Processing, 2021, 15(4): 1029-1042.

- [4] Snyder D, Garcia-Romero D, Sell G, et al. X-vectors: Robust DNN embeddings for speaker recognition. IEEE Transactions on Audio, Speech, and Language Processing, 2018, 26(7): 110-119.
- [5] Chung J S, Nagrani A, Zisserman A. VoxCeleb2: Deep speaker recognition. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2021, 13(5): 532-541.
- [6] Villalba J. Advanced speaker recognition using deep neural networks. Carnegie Mellon University, 2020.
- [7] Hajibabaei M, Dai D. Unified hypersphere embedding for speaker recognition. IEEE Transactions on Audio, Speech, and Language Processing, 2019, 45(3): 321-330.
- [8] Liu X, Ji Y, Liu C. Voiceprint recognition based on LSTM neural network . Computer Science, 2021, 48(S2), 270-274.
- [9] Wang H P. Speaker recognition based on deep bidirectional LSTM network. Computer Engineering and Design, 2020, 41(06): 1768-1772.
- [10] Zhao H, Yue L, Wang W, et al. Research on end-to-end voiceprint recognition model based on convolutional neural network. Journal of Web Engineering, 2021, 20(5): 1573-1586.
- [11] Yan H, Dong Y, Wang P, et al. Research on voiceprint recognition based on CNN-LSTM network. Computer Application and Software, 2019, 36(04): 166-170.
- [12] Guo D, Zhou Q. Land-air call voiceprint recognition based on residual neural network. Modern Computer, 2020(07): 9-13.
- [13] Liu Y, Liang H, Liu G, et al. Voiceprint recognition method based on ResNet-LSTM. Computer System Applications, 2021, 30(06): 215-219.