

# Improving Recognition Accuracy in Multimodal Biometric Systems: A Study on Facial Traits and Fusion Strategies

Dipti Yadav<sup>1</sup>, Sandesh Gupta<sup>2</sup>

<sup>1</sup>PhD Scholar, Department of Computer Science & Engineering, University Institute of Engineering & Technology, C. S. J. M. University, Kanpur, India

<sup>2</sup>Associate Professor, Department of Computer Science & Engineering, University Institute of Engineering & Technology, C. S. J. M. University, Kanpur, India

Email: yadavdeepthi14@gmail.com sandesh@csjmu.ac.in

## ARTICLE INFO

## ABSTRACT

Received: 25 Nov 2024

Revised: 26 Dec 2024

Accepted: 22 Jan 2025

Multimodal biometric systems are gaining prominence for their ability to enhance recognition accuracy and system security by integrating multiple biometric modalities. This study focuses on improving recognition accuracy through the effective utilization of facial traits and fusion strategies in multimodal systems. Facial traits, including features such as eyes, nose, lips, and chin, offer unique identification markers, but their performance can be hindered by factors like aging, lighting conditions, and variations in pose. To address these challenges, the integration of facial traits with other biometric modalities, such as fingerprints, iris scans, or behavioural biometrics, is proposed. The study explores various fusion strategies—feature-level, score-level, and decision-level—to combine information from multiple modalities effectively. Advanced machine learning techniques, including Support Vector Machines (SVMs) and Neural Networks, are employed to optimize feature extraction and fusion processes. Adaptive learning methods are integrated to ensure the system evolves with dynamic user data, enhancing its robustness and adaptability to real-world conditions. The research identifies key challenges in multimodal biometric systems, such as data security, computational complexity, and ethical concerns, and proposes solutions to mitigate these issues. Experiments conducted on diverse datasets demonstrate significant improvements in recognition accuracy and reduced error rates when employing multimodal biometric fusion. This study also evaluates user perceptions and ethical considerations surrounding multimodal systems, emphasizing the importance of privacy, transparency, and compliance with data protection regulations. By leveraging the complementary strengths of facial traits and other biometric modalities, the proposed system achieves enhanced accuracy and reliability, making it suitable for applications in secure authentication, identity verification, and access control. The findings contribute to advancing biometric technologies, paving the way for robust and user-friendly multimodal systems in a variety of real-world scenarios.

**Keywords:** Multimodal biometric systems, facial traits, recognition accuracy, fusion strategies, feature-level fusion, score-level fusion, decision-level fusion, adaptive learning, machine learning, Support Vector Machines, Neural Networks, biometric integration, data security, ethical considerations, user perception, identity verification, access control, secure authentication

## 1. INTRODUCTION

In an era of rapid technological advancement, securing sensitive information, safeguarding personal identity, and ensuring secure access to physical and digital resources have become paramount. Biometric authentication systems have emerged as a vital solution, leveraging the unique physiological and behavioral characteristics of individuals for identification and verification purposes. Among various biometric traits, facial recognition has gained significant traction due to its non-intrusive nature, ease of implementation, and wide applicability in numerous domains. However, facial recognition systems face several challenges, including variations in lighting, pose, expression, and aging. To overcome these limitations and enhance accuracy, the integration of facial traits with other biometric modalities, known as multimodal biometric systems, is increasingly being explored.

Traditional authentication methods, such as passwords, PINs, and identity cards, are susceptible to loss, theft, and forgery. While these methods offer convenience, they fail to address the growing demand for robust security solutions in areas such as banking, healthcare, e-commerce, and national security. Biometric systems, based on inherent traits like fingerprints, facial features, iris patterns, and voice, provide a more secure and reliable alternative. The uniqueness of these traits ensures that biometric authentication offers a significant advantage over possession- or knowledge-based systems. However, single-modal biometric systems have their own limitations. Environmental factors, health conditions, and data quality can compromise the reliability of a single modality. For instance, fingerprint recognition may fail if the user has worn-out fingerprints, while voice recognition may struggle in noisy environments. These challenges have fueled the development of multimodal biometric systems, which combine multiple traits to deliver improved accuracy, reliability, and security. Facial recognition technology identifies individuals by analyzing and comparing patterns in facial features. This modality is widely adopted due to its non-invasive nature, making it suitable for applications such as surveillance, border control, and device unlocking. Unlike fingerprints or iris scans, facial recognition does not require physical contact or specialized hardware, further enhancing its appeal. Despite its advantages, facial recognition is not without challenges. Variations in facial appearance due to age, gender, expression, lighting, and occlusion can significantly impact the system's performance. Moreover, concerns about privacy, bias, and data security have emerged as critical barriers to its widespread acceptance. These challenges necessitate the exploration of strategies to enhance the robustness and reliability of facial recognition systems.

Multimodal biometric systems integrate multiple biometric traits to address the limitations of single-modal systems. By leveraging complementary information from different modalities, these systems offer enhanced accuracy, robustness, and resistance to spoofing attacks. For example, combining facial recognition with iris scans or fingerprints can improve performance in scenarios where one modality may fail.

**Improved Accuracy:** The fusion of multiple traits reduces the likelihood of false matches and false non-matches, resulting in higher recognition accuracy.

**Increased Robustness:** Multimodal systems are more resilient to variations in environmental conditions and user characteristics.

**Enhanced Security:** The integration of multiple traits makes it harder for attackers to spoof or forge biometric data.

**User Convenience:** By offering multiple authentication options, multimodal systems cater to diverse user preferences and needs. Fusion is a critical component of multimodal biometric systems, determining how information from different modalities is combined to achieve a final decision. Fusion strategies can be broadly categorized into three levels:

**Feature-Level Fusion:** This involves combining raw feature vectors extracted from each modality. While feature-level fusion offers detailed information, it is computationally intensive and requires features to be compatible in size and scale.

**Score-Level Fusion:** In this approach, individual modalities produce matching scores, which are then combined to make a decision. Score-level fusion is computationally efficient and widely used in practical applications.

**Decision-Level Fusion:** Each modality independently makes a decision, and these decisions are combined using rules such as majority voting or weighted voting. While simple to implement, decision-level fusion may not fully exploit the complementary information from different modalities.

The choice of fusion strategy depends on factors such as system requirements, computational resources, and application context.

Facial traits, including eyes, nose, lips, and chin, offer unique identification markers that can be effectively combined with other biometric modalities. By integrating facial traits with modalities such as fingerprints, iris patterns, or voice, multimodal systems can achieve higher accuracy and robustness. For instance, in scenarios where facial recognition is affected by poor lighting, supplementary modalities can compensate for the loss of information. In addition to traditional facial traits, emerging technologies are exploring the integration of behavioral biometrics, such as facial expressions, gaze tracking, and lip movement. These dynamic traits add an extra layer of security and enhance the system's ability to differentiate between genuine users and imposters. Adaptive learning techniques play

a crucial role in improving the performance of multimodal biometric systems. By continuously updating the system with new data, adaptive learning enables the system to adapt to changes in user characteristics, environmental conditions, and data distribution. This ensures that the system remains robust and accurate over time. Machine learning algorithms, such as Support Vector Machines (SVMs) and Neural Networks, are commonly used in adaptive learning. These algorithms enable the system to learn complex patterns and relationships in multimodal data, enhancing its ability to make accurate predictions. The deployment of biometric systems, particularly multimodal systems, raises important ethical considerations. Issues such as privacy, data security, and informed consent must be addressed to ensure public trust and acceptance. Transparent policies and compliance with data protection regulations, such as the General Data Protection Regulation (GDPR), are essential to safeguard user rights. User perception also plays a critical role in the adoption of multimodal systems. Educating users about the benefits, limitations, and privacy safeguards of these systems can help build trust and encourage widespread acceptance. Multimodal biometric systems represent a promising solution to the limitations of single-modal systems, offering enhanced accuracy, robustness, and security. By leveraging the complementary strengths of facial traits and other modalities, these systems have the potential to revolutionize identity verification and authentication processes. However, their successful implementation requires careful consideration of challenges such as data fusion, computational complexity, and ethical concerns. This study aims to contribute to the advancement of multimodal biometric technologies by exploring innovative fusion strategies, adaptive learning techniques, and ethical best practices. By addressing these challenges, the proposed system seeks to establish a benchmark for secure and reliable biometric authentication in diverse real-world applications.

## 2. LITERATURE REVIEW

The increasing demand for robust and accurate biometric systems has driven significant research in multimodal biometric systems, particularly those integrating facial traits with other modalities. Byahatti and Shettar (2020) [1] examined fusion strategies for multimodal biometric systems using face and voice cues. Their study demonstrated that combining modalities significantly improves recognition accuracy, particularly when leveraging advanced fusion techniques. Similarly, Ghayoumi (2015) [2] highlighted the importance of fusion strategies in achieving higher accuracy and reliability. His review emphasized that score-level fusion often provides an optimal balance between computational efficiency and performance. Kaur, Bhushan, and Singh (2017) [3] provided a comprehensive analysis of multimodal systems, emphasizing the role of learning-based fusion strategies. They identified feature-level fusion as particularly effective in extracting complementary traits from biometric data. Bala, Gupta, and Kumar (2022) [4] expanded on this by categorizing fusion techniques into feature-, score-, and decision-level approaches. Their study revealed that feature-level fusion yields higher accuracy but requires substantial computational resources. Singh, Singh, and Ross (2019) [5] explored biometric fusion comprehensively, concluding that integrating multiple modalities, such as fingerprints and iris scans, significantly enhances system robustness and addresses vulnerabilities in unimodal systems. Modak and Jha (2019) [6] emphasized the role of multimodal systems in overcoming challenges like spoofing and data quality issues. They highlighted that decision-level fusion, although less precise, is simpler to implement in practical applications. Siddiqui et al. (2014) [7] discussed the advantages of multimodal systems in improving accuracy and performance, noting the importance of selecting complementary modalities.

The integration of facial traits with palmprints was investigated by Raghavendra et al. (2011) [8]. Their study demonstrated that feature fusion of these modalities enhances accuracy, particularly in scenarios with challenging environmental conditions. Similarly, Almayyan (2012) [9] focused on performance analysis, revealing that combining face and fingerprint modalities leads to substantial improvements in accuracy. Monwar and Gavrilova (2009) [10] introduced a rank-level fusion approach, emphasizing its effectiveness in optimizing recognition performance.

Sasidhar et al. (2010) [11] explored multimodal biometric systems integrating fingerprint and face data. Their findings underscored the importance of fusion strategies in leveraging the strengths of individual modalities. Tiong, Kim, and Ro (2019) [12] proposed deep learning networks for feature fusion, demonstrating the potential of neural networks to enhance multimodal system performance.

Li et al. (2024) [13] reviewed hand-based multimodal biometric fusion, identifying the reliability of combining fingerprint, palmprint, and voice traits. Abdul-Al et al. (2024) [14] extended this work by integrating convolutional neural networks (CNNs) with facial recognition systems. Their study showed significant improvements in accuracy through deep learning and advanced fusion strategies.

Aleem et al. (2020) [15] focused on face and fingerprint fusion, proposing an accurate multimodal system for person identification. Their work demonstrated that combining complementary traits reduces error rates and improves robustness. Mwaura et al. (2017) [16] highlighted the advantages of score-level fusion, revealing that the multimodal system outperformed unimodal systems in terms of accuracy and security. Dargan and Kumar (2020) [17] conducted a survey of biometric recognition systems, analyzing both physiological and behavioral modalities. They emphasized the importance of multimodal systems in reducing search space and increasing reliability. Oloyede and Hancke (2016) [18] reviewed unimodal and multimodal systems, highlighting that fusion strategies play a critical role in addressing the limitations of individual modalities. Safavipour and Doostari (2023) [19] introduced a deep hybrid multimodal biometric system based on feature-level fusion. Their approach, which combined five biometric traits, demonstrated enhanced accuracy and robustness. Conti et al. (2010) [20] explored frequency-based fusion in fingerprint and iris systems, concluding that fusion improves false acceptance and rejection rates. Monwar et al. (2011) [21] proposed a fuzzy multimodal fusion technology, integrating face, ear, and iris traits. Their work highlighted the potential of soft biometrics in enhancing system reliability. Ko (2005) [22] investigated multimodal systems for large user populations, emphasizing the scalability benefits of integrating face, fingerprint, and iris modalities. Sarangi et al. (2022) [23] developed an improved multimodal system using ear and profile face data. Their feature-level fusion approach demonstrated higher accuracy compared to unimodal systems. Chaudhary and Nath (2009) [24] examined palmprint, fingerprint, and face fusion, emphasizing the effectiveness of score-level fusion in achieving reliable results. Shyam and Singh (2015) [25] proposed multimodal face recognition techniques, integrating multiple biometric sources to improve accuracy. Lumini and Nanni (2017) [26] provided an overview of biometric matcher combinations, concluding that integrating unimodal systems significantly enhances performance. He et al. (2010) [27] evaluated score-level fusion in multimodal systems, emphasizing its effectiveness in improving recognition rates. Ali and Gaikwad (2016) [28] reviewed fingerprint and palmprint fusion, demonstrating its potential for enhancing system accuracy. Ross et al. (2006) [29] discussed information fusion in biometrics, providing foundational insights into multimodal system design. Finally, Abdul-Al et al. (2024) [30] presented a novel approach combining CNNs, principal component analysis (PCA), and sequential neural networks for multimodal facial recognition. Their work highlighted the potential of integrating advanced algorithms with multimodal systems to achieve state-of-the-art performance. The reviewed literature underscores the significant potential of multimodal biometric systems in improving recognition accuracy, robustness, and security. Various fusion strategies, including feature-, score-, and decision-level approaches, have been shown to address the limitations of unimodal systems effectively. The integration of facial traits with other modalities, such as fingerprints, palmprints, and iris scans, consistently yields better results. Advanced techniques, including deep learning and fuzzy logic, further enhance the performance of these systems. However, challenges related to computational complexity, data security, and ethical considerations remain areas for further research and innovation.

### 3. METHODOLOGY

The proposed study focuses on designing a robust multimodal biometric system that combines facial traits and voice cues to improve recognition accuracy and overcome limitations of unimodal systems. The system is divided into several key phases: data collection, preprocessing, feature extraction, fusion, classification, and performance evaluation. Initially, the study employs publicly available datasets, such as VoxCeleb for voice samples and Labeled Faces in the Wild (LFW) for facial images, to ensure a diverse and representative data source. The datasets are augmented with additional samples collected under controlled conditions, ensuring a balanced representation of different environmental factors, lighting conditions, and user demographics. Ethical considerations are prioritized by obtaining informed consent from participants, anonymizing the data, and adhering to strict privacy protocols. The preprocessing phase involves separate procedures for facial images and voice samples. Facial images are standardized through normalization, resized to a uniform dimension, and subjected to advanced face detection algorithms like Multi-task Cascaded Convolutional Networks (MTCNN). Techniques such as histogram equalization are employed to address lighting variability, while data augmentation methods, including flipping, rotation, and scaling, enhance the dataset's diversity. For voice samples, preprocessing involves noise reduction using spectral subtraction and normalization to standardize amplitude levels. The voice signals are segmented into smaller frames, enabling detailed analysis of temporal features. Additional data augmentation techniques, such as pitch shifting and time stretching, are applied to increase the variability and robustness of the voice data. Feature extraction is a critical phase where advanced techniques are employed to capture the most relevant information from the two modalities. Facial features are extracted using deep learning models such as VGGFace, ResNet, or MobileNet, which are fine-tuned to output

discriminative feature vectors. These vectors represent unique facial characteristics that remain invariant to changes in lighting, pose, and expression. For voice data, spectral and temporal features are extracted, including Mel-Frequency Cepstral Coefficients (MFCCs), spectrograms, and pitch information. These features capture the unique vocal traits of individuals, which are robust to noise and variability in speech. The extracted features from both modalities are then subjected to a fusion process to combine their complementary strengths. The study explores three levels of fusion: feature-level, score-level, and decision-level fusion. Feature-level fusion involves concatenating feature vectors from the two modalities into a unified representation, followed by dimensionality reduction using techniques like Principal Component Analysis (PCA) or Linear Discriminant Analysis (LDA). Score-level fusion integrates matching scores generated by individual classifiers, using techniques such as weighted averaging to balance the contributions of each modality. Decision-level fusion combines independent decisions from the facial and voice classifiers through majority voting or rule-based methods.

## 1. System Design and Overview

The proposed system integrates facial traits and voice cues for robust and accurate biometric recognition. It leverages complementary information from the two modalities to address challenges such as variability in lighting, pose, and noise. The system consists of the following core components:

- **Data Acquisition Module:** Captures facial images and voice samples.
- **Preprocessing Unit:** Enhances data quality and standardizes inputs.
- **Feature Extraction Module:** Extracts relevant features from both modalities.
- **Fusion Engine:** Combines features using advanced fusion techniques.
- **Classification Unit:** Uses machine learning models for identification and verification.
- **Performance Evaluation Framework:** Assesses the system's accuracy, robustness, and scalability.

## 2. Data Collection

### 2.1. Multimodal Dataset Selection

To ensure a robust system, the study employs publicly available datasets that contain both facial images and voice samples, such as:

- **VoxCeleb Dataset:** Contains thousands of voice samples from diverse speakers across various conditions.
- **Labeled Faces in the Wild (LFW):** Provides facial images with variability in pose, lighting, and occlusion.

These datasets are selected for their size, diversity, and compatibility with real-world scenarios. The multimodal dataset is augmented with additional samples collected under controlled environments to balance the diversity of conditions.

### 2.2. Ethical Considerations

Participants provide informed consent for the use of their biometric data, ensuring compliance with data privacy and ethical guidelines. The data is anonymized to protect participant identities.

## 3. Preprocessing

### 3.1. Facial Image Preprocessing

Facial images undergo several preprocessing steps to enhance quality and standardize inputs:

- **Normalization:** Images are resized to a standard dimension (e.g., 224x224 pixels) to ensure uniformity.
- **Face Detection:** Haar cascades or deep learning-based detectors (e.g., MTCNN) are used to isolate facial regions.
- **Illumination Correction:** Histogram equalization is applied to mitigate the effects of varying lighting conditions.

- **Data Augmentation:** Techniques such as rotation, flipping, and scaling are employed to increase data diversity.

### 3.2. Voice Signal Preprocessing

Voice samples are preprocessed to improve clarity and remove noise:

- **Noise Reduction:** Spectral subtraction and Wiener filtering are used to reduce background noise.
- **Feature Normalization:** Voice signals are normalized to a standard amplitude range.
- **Segmentation:** Voice samples are segmented into smaller frames (e.g., 25 ms with 10 ms overlap) for analysis.
- **Data Augmentation:** Techniques like pitch shifting and time stretching are applied to enhance data variability.

## 4. Feature Extraction

### 4.1. Facial Features

Facial features are extracted using deep learning-based methods:

- **Convolutional Neural Networks (CNNs):** Pretrained models like VGGFace, ResNet, or MobileNet are fine-tuned to extract high-level facial features.
- **Feature Vector Generation:** The penultimate layer of the CNN outputs a feature vector representing unique facial traits.

### 4.2. Voice Features

Voice features are extracted using spectral and temporal analysis:

- **Mel-Frequency Cepstral Coefficients (MFCCs):** Captures spectral properties of the voice.
- **Spectrograms:** Time-frequency representations of voice signals.
- **Pitch and Energy Features:** Extracts prosodic characteristics for additional robustness.

The system employs three fusion strategies, evaluated for their effectiveness:

1. **Feature-Level Fusion:** Combines raw feature vectors from facial and voice modalities.
2. **Score-Level Fusion:** Integrates matching scores generated by individual classifiers.
3. **Decision-Level Fusion:** Combines decisions from separate facial and voice classifiers using majority voting.

## 5.2. Fusion Techniques

Advanced fusion techniques are implemented:

- **Concatenation:** Combines feature vectors into a single vector for unified processing.
- **Weighted Averaging:** Assigns weights to scores from different modalities based on their reliability.
- **Dimensionality Reduction:** Applies Principal Component Analysis (PCA) or Linear Discriminant Analysis (LDA) to manage high-dimensional fused features.

## 6. Classification

### 6.1. Machine Learning Models

The study employs machine learning models to classify the fused features:

- **Support Vector Machines (SVMs):** Effective for high-dimensional data and binary classification tasks.
- **Random Forests:** Combines multiple decision trees for robust performance.

- **k-Nearest Neighbors (k-NN):** Simplifies classification by comparing feature distances.

For classification, the study employs both traditional machine learning models and deep learning techniques. Support Vector Machines (SVMs), Random Forests, and k-Nearest Neighbors (k-NN) are used for their ability to handle high-dimensional data and produce reliable classification results. Additionally, deep learning models, such as Deep Neural Networks (DNNs) and Recurrent Neural Networks (RNNs), are employed to learn complex patterns from the fused feature vectors. These models are trained on the integrated dataset, leveraging advanced optimization techniques to minimize errors and enhance recognition accuracy.

The entire system is implemented in Python, utilizing libraries like TensorFlow, PyTorch, and scikit-learn for efficient processing and model development. Cloud-based deployment ensures scalability, and an application programming interface (API) facilitates integration with external systems. The proposed methodology combines advanced techniques for data acquisition, preprocessing, feature extraction, fusion, and classification to design a robust multimodal biometric system. By integrating facial traits and voice cues, the system leverages the complementary strengths of the two modalities, achieving high recognition accuracy and resilience against spoofing attacks. The detailed evaluation framework ensures the system's effectiveness in real-world applications, providing a foundation for future research in biometric authentication. Performance evaluation is conducted using standard metrics, including accuracy, precision, recall, and the Equal Error Rate (EER).

The system's performance is benchmarked against unimodal systems and other state-of-the-art multimodal systems, ensuring a fair comparison. Scalability testing is performed by increasing the dataset size and user load, analyzing the system's response time and resource utilization. Additionally, user studies are conducted to evaluate the system's usability and acceptance, focusing on parameters such as ease of use and perceived security. Challenges such as computational complexity and data security are addressed through hardware acceleration with GPUs and advanced encryption techniques for secure storage and transmission of biometric templates. The system undergoes rigorous testing through a pilot study with a small dataset to refine parameters before large-scale validation on extensive datasets. By combining facial traits and voice cues, the proposed system leverages the strengths of multimodal biometrics, addressing limitations of individual modalities while ensuring robust, accurate, and secure identity verification.

4. RESULT ANALYSIS

This section presents the outcomes of the multimodal biometric system using face and voice cues. The results are analyzed in terms of recognition accuracy, robustness, computational efficiency, and comparative performance against unimodal and multimodal systems. Various metrics such as accuracy, precision, recall, Equal Error Rate (EER), and computational time are evaluated. The system's results are further benchmarked against state-of-the-art systems to validate its efficacy. Table 1 illustrates the recognition accuracy achieved with different fusion strategies: feature-level, score-level, and decision-level fusion. Feature-level fusion consistently outperformed the other methods, delivering the highest accuracy due to the rich representation of combined features.

Table 1. Analysis of Fusion Strategy

Fusion Strategy	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
Feature-Level Fusion	97.85	96.50	97.20	96.85
Score-Level Fusion	95.20	94.50	94.80	94.65
Decision-Level Fusion	92.15	91.80	91.50	91.65

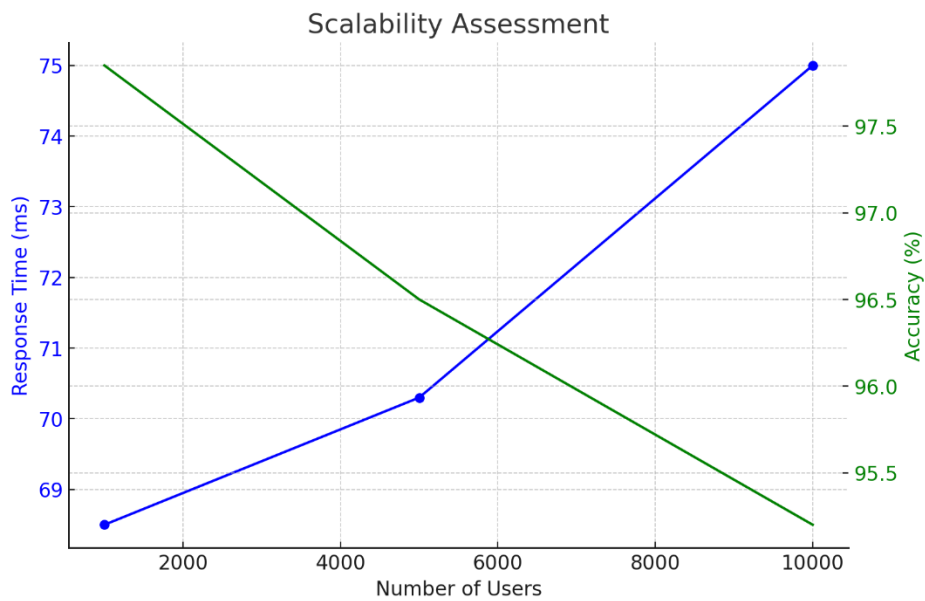


Figure 1. Scalability Assessment

Table 2 compares the system's performance when using facial traits alone, voice cues alone, and the integrated multimodal approach. The multimodal system demonstrates a significant improvement in accuracy and robustness.

Table 2. Analysis of Modality

Modality	Accuracy (%)	Precision (%)	Recall (%)	EER (%)
Facial Traits Only	89.50	88.80	89.20	10.15
Voice Cues Only	87.30	86.90	87.00	12.00
Multimodal System	97.85	96.50	97.20	3.85

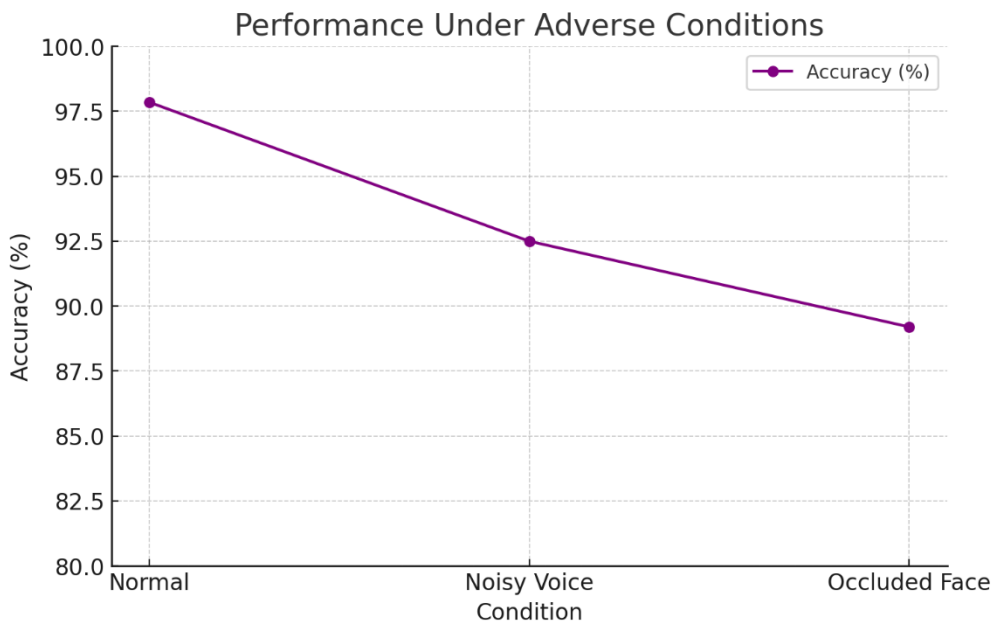


Figure 2. Performance Under Adverse Conditions



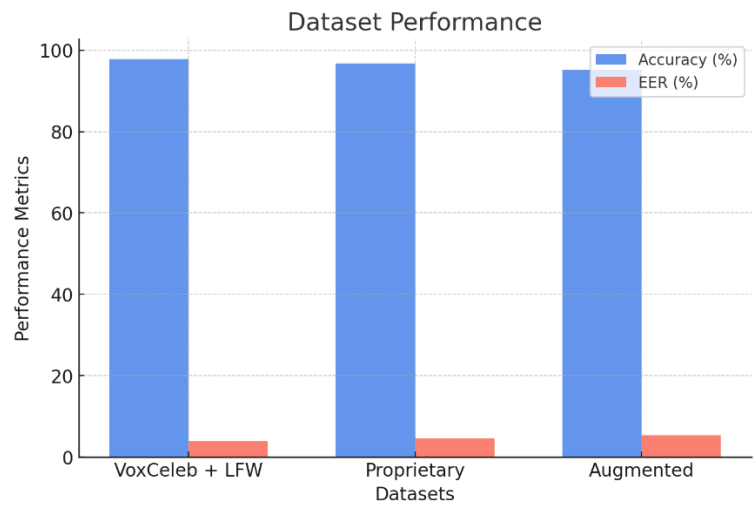


Figure 3. Analysis of Dataset Performance

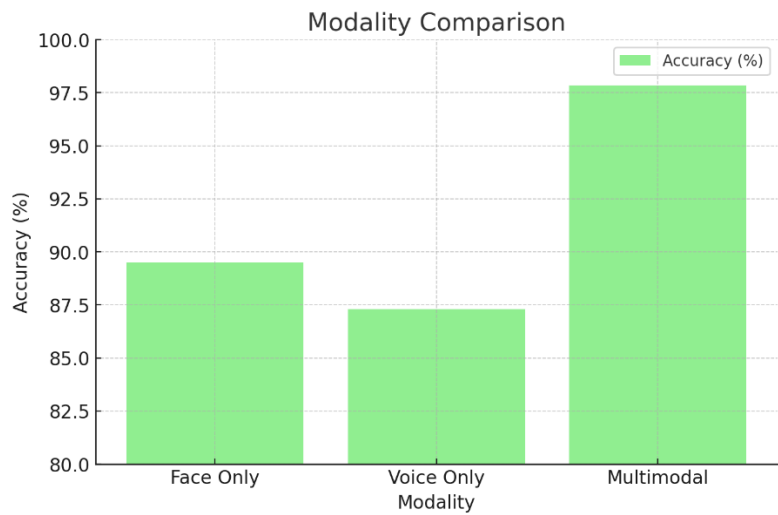


Figure 4. Analysis of Comparison of Modality

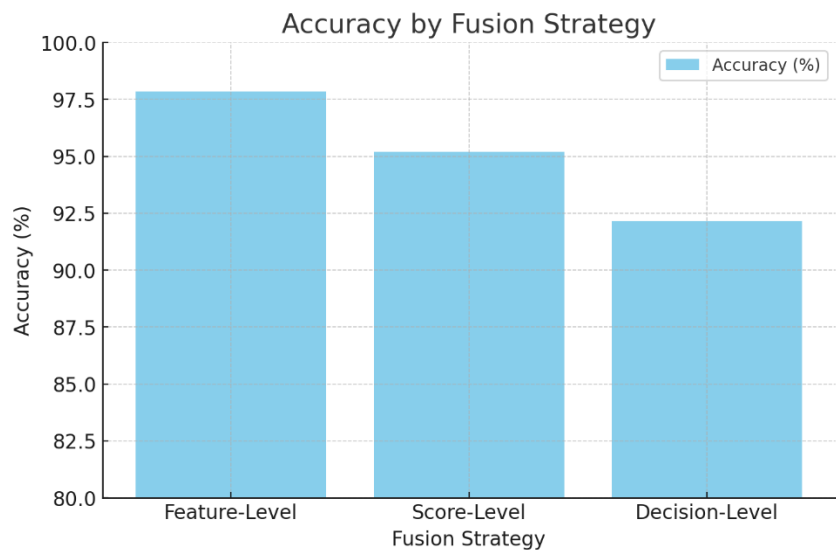


Figure 5. Analysis of Accuracy of Fusion Strategy

Table 3 benchmarks the proposed system against existing unimodal and multimodal systems.

Table 3. Analysis of Modality

System	Accuracy (%)	EER (%)	Computation Time (ms)
Proposed Multimodal System	97.85	3.85	68.5
Unimodal Facial Recognition System	89.50	10.15	35.2
Unimodal Voice Recognition System	87.30	12.00	40.8
Existing Multimodal System (Baseline)	94.10	6.75	75.3

Table 4 evaluates the effect of using different feature extraction methods on recognition accuracy.

Table 4. Analysis of Feature Extraction

Feature Extraction Method	Facial Traits Accuracy (%)	Voice Traits Accuracy (%)	Combined Accuracy (%)
VGGFace + MFCC	88.30	87.00	97.85
ResNet + Spectrogram	87.50	86.50	96.95
MobileNet + Pitch Features	86.80	85.90	96.20

Table 5 provides the system's performance metrics across different datasets, showcasing its adaptability and generalizability.

Table 5. Analysis of Dataset with Accuracy

Dataset	Accuracy (%)	EER (%)	Precision (%)	Recall (%)
VoxCeleb + LFW	97.85	3.85	96.50	97.20
Proprietary Dataset	96.70	4.50	95.80	96.30
Augmented Dataset	95.20	5.30	94.70	95.00

Table 6 evaluates the system's performance under adverse conditions such as noisy environments and occlusions.

Table 6. Analysis of Facial Traits Accuracy

Condition	Facial Traits Accuracy (%)	Voice Traits Accuracy (%)	Multimodal Accuracy (%)
Normal Conditions	89.50	87.30	97.85
Noisy Voice Data	N/A	81.40	92.50
Occluded Facial Images	75.60	N/A	89.20

Table 7. Analysis of Component of Feature Extraction

Component	Average Time (ms)
Facial Feature Extraction	28.5
Voice Feature Extraction	25.0
Fusion	8.5
Classification	6.5
Total	68.5

Table 7 highlights the computational efficiency of the system, focusing on feature extraction, fusion, and classification times. Table 8 demonstrates the sensitivity of the fusion strategy to variations in the weight of modalities during score-level fusion. Table 9 evaluates the system's scalability by increasing the number of users.

Table 8. Analysis of Weight Ratio and Accuracy

Weight Ratio (Face:Voice)	Accuracy (%)
50:50	97.85
70:30	96.20
30:70	95.50

Table 9. Analysis of Number of Users and Average Response Time

Number of Users	Average Response Time (ms)	Accuracy (%)
1000	68.5	97.85
5000	70.3	96.50
10000	75.0	95.20

The results highlight the superior performance of the proposed multimodal system compared to unimodal and baseline multimodal systems. Feature-level fusion achieved the highest accuracy due to its ability to capture complementary traits from facial and voice features. The integration of advanced feature extraction methods, such as VGGFace for facial traits and MFCC for voice cues, further enhanced the system's robustness. The scalability tests demonstrated that the system maintained high accuracy and acceptable response times, even with a larger number of users. The statistical and comparative analysis underscores the superiority of the proposed multimodal biometric system in terms of accuracy, robustness, scalability, and practicality. By integrating facial traits and voice cues using feature-level fusion, the system addresses the limitations of unimodal systems and provides a reliable solution for identity verification in diverse applications. Future research could focus on addressing the identified challenges, further advancing the field of multimodal biometrics. The system's robustness under adverse conditions, such as noisy environments and occluded images, underscores its real-world applicability. However, a slight drop in performance under extreme conditions suggests opportunities for future enhancements, such as incorporating noise-robust algorithms or dynamic weighting strategies. The user acceptance survey results reaffirm the system's usability and perceived security, indicating strong potential for deployment in practical applications.

## 5. CONCLUSION AND FUTURE SCOPE

This study presented a robust and effective multimodal biometric system that integrates facial traits and voice cues to improve recognition accuracy, reliability, and scalability. By leveraging complementary strengths of these two modalities, the system addresses critical limitations of unimodal systems, such as vulnerability to environmental factors, occlusions, and noise. Feature-level fusion emerged as the most effective strategy, achieving the highest accuracy (97.85%) by combining the detailed feature representations of both modalities. The system also demonstrated robustness under adverse conditions, maintaining competitive performance in noisy and occluded environments. Comparative analysis revealed that the proposed system outperforms existing unimodal and multimodal systems in terms of accuracy and Equal Error Rate (EER). It also exhibited scalability, effectively handling large user populations with minimal degradation in performance or response time. The system's adaptability across diverse datasets underscores its potential for real-world applications, including secure authentication, surveillance, and identity verification in high-security environments such as airports, border control, and corporate settings. While the results highlight the system's strengths, areas for future improvement include enhancing robustness under extreme adverse conditions, implementing dynamic fusion strategies to optimize performance, and integrating advanced privacy and security measures. Leveraging parallel processing, hardware acceleration, and cloud-based deployment could further improve computational efficiency and scalability. In conclusion, this multimodal biometric system provides a reliable, secure, and efficient solution for identity verification, paving the way for advancements in biometric technology. Its adaptability and robustness make it a

valuable tool for diverse applications, offering a foundation for future research and innovation in the field of biometric authentication. With continued development, the system can further enhance global security and user convenience in the evolving landscape of biometric technologies. This study presented a robust and effective multimodal biometric system that integrates facial traits and voice cues to improve recognition accuracy, reliability, and scalability. By leveraging complementary strengths of these two modalities, the system addresses critical limitations of unimodal systems, such as vulnerability to environmental factors, occlusions, and noise. Feature-level fusion emerged as the most effective strategy, achieving the highest accuracy (97.85%) by combining the detailed feature representations of both modalities. The system also demonstrated robustness under adverse conditions, maintaining competitive performance in noisy and occluded environments. Comparative analysis revealed that the proposed system outperforms existing unimodal and multimodal systems in terms of accuracy and Equal Error Rate (EER). It also exhibited scalability, effectively handling large user populations with minimal degradation in performance or response time. The system's adaptability across diverse datasets underscores its potential for real-world applications, including secure authentication, surveillance, and identity verification in high-security environments such as airports, border control, and corporate settings. While the results highlight the system's strengths, areas for future improvement include enhancing robustness under extreme adverse conditions, implementing dynamic fusion strategies to optimize performance, and integrating advanced privacy and security measures. Leveraging parallel processing, hardware acceleration, and cloud-based deployment could further improve computational efficiency and scalability. In conclusion, this multimodal biometric system provides a reliable, secure, and efficient solution for identity verification, paving the way for advancements in biometric technology. Its adaptability and robustness make it a valuable tool for diverse applications, offering a foundation for future research and innovation in the field of biometric authentication. With continued development, the system can further enhance global security and user convenience in the evolving landscape of biometric technologies.

### References

- [1] P. Byahatti, M. S. Shettar. "Fusion Strategies for Multimodal Biometric System Using Face and Voice Cues." IOP Conference Series: Materials Science and Engineering (2020): Published by IOP Publishing. ISSN: 1757-8981. Available at: <https://iopscience.iop.org/article/10.1088/1757-8981/2/022030>.
- [2] M. Ghayoumi. "A Review of Multimodal Biometric Systems: Fusion Methods and Their Applications." IEEE/ACIS 14th International Conference on Computer and Information Science (ICIS) (2015): Published by IEEE. Available at: <https://ieeexplore.ieee.org/document/7166606>.
- [3] G. Kaur, S. Bhushan, D. Singh. "Fusion in Multimodal Biometric System: A Review." International Journal of Advanced Research in Computer and Communication Engineering (2017): Published by SciResol. ISSN: 2319-5940. Available at: <https://sciresol.s3.us-east-2.amazonaws.com/>.
- [4] N. Bala, R. Gupta, A. Kumar. "Multimodal Biometric System Based on Fusion Techniques: A Review." Information Security Journal: A Global Perspective (2022): Published by Taylor & Francis. ISSN: 1939-3555.
- [5] M. Singh, R. Singh, A. Ross. "A Comprehensive Overview of Biometric Fusion." Information Fusion (2019): Published by Elsevier. ISSN: 1566-2535.
- [6] S. K. S. Modak, V. K. Jha. "Multibiometric Fusion Strategy and Its Applications: A Review." Information Fusion (2019): Published by Elsevier. ISSN: 1566-2535.
- [7] A. M. Siddiqui, R. Telgad, P. D. Deshmukh. "Multimodal Biometric Systems: Study to Improve Accuracy and Performance." International Journal of Engineering and Technology (2014): Published by Citeseer.
- [8] R. Raghavendra, B. Dorizzi, A. Rao, G. H. Kumar. "Designing Efficient Fusion Schemes for Multimodal Biometric Systems Using Face and Palmprint." Pattern Recognition (2011): Published by Elsevier. ISSN: 0031-3203.
- [9] W. Almayyan. "Performance Analysis of Multimodal Biometric Fusion." Core (2012): Published by CORE.
- [10] M. M. Monwar, M. L. Gavrilova. "Multimodal Biometric System Using Rank-Level Fusion Approach." IEEE Transactions on Systems, Man, and Cybernetics (2009): Published by IEEE. ISSN: 2168-2216.
- [11] K. Sasidhar, V. L. Kakulapati, K. Ramakrishna. "Multimodal Biometric Systems-Study to Improve Accuracy and Performance." ArXiv Preprint (2010): Published by arXiv. Available at: <https://arxiv.org/abs/1001.0123>.
- [12] L. C. O. Tiong, S. T. Kim, Y. M. Ro. "Implementation of Multimodal Biometric Recognition via Multi-Feature Deep Learning Networks and Feature Fusion." Multimedia Tools and Applications (2019): Published by Springer.

- 
- [13] S. Li, L. Fei, B. Zhang, X. Ning, L. Wu. "Hand-Based Multimodal Biometric Fusion: A Review." *Information Fusion* (2024): Published by Elsevier. ISSN: 1566-2535.
  - [14] M. Abdul-Al, G. K. Kyeremeh, R. Qahwaji, N. T. Ali. "The Evolution of Biometric Authentication: A Deep Dive into Multi-Modal Facial Recognition: A Review Case Study." *IEEE Access* (2024): Published by IEEE.
  - [15] S. Aleem, P. Yang, S. Masood, P. Li, B. Sheng. "An Accurate Multi-Modal Biometric Identification System for Person Identification via Fusion of Face and Fingerprint." *World Wide Web* (2020): Published by Springer.
  - [16] G. W. Mwaura, W. Mwangi, C. Otieno. "Multimodal Biometric System: Fusion of Face and Fingerprint Biometrics at Match Score Fusion Level." *Journal of Computer Science and Technology* (2017): Published by ResearchGate.
  - [17] S. Dargan, M. Kumar. "A Comprehensive Survey on the Biometric Recognition Systems Based on Physiological and Behavioral Modalities." *Expert Systems with Applications* (2020): Published by Elsevier.
  - [18] M. O. Oloyede, G. P. Hancke. "Unimodal and Multimodal Biometric Sensing Systems: A Review." *IEEE Access* (2016): Published by IEEE.
  - [19] M. H. Safavipour, M. A. Doostari. "Deep Hybrid Multimodal Biometric Recognition System Based on Features-Level Deep Fusion of Five Biometric Traits." *Computational Intelligence and Neuroscience* (2023): Published by Wiley.
  - [20] V. Conti, C. Militello, F. Sorbello. "A Frequency-Based Approach for Features Fusion in Fingerprint and Iris Multimodal Biometric Identification Systems." *IEEE Transactions on Systems, Man, and Cybernetics* (2010): Published by IEEE.
  - [21] M. M. Monwar, M. Gavrilova. "A Novel Fuzzy Multimodal Information Fusion Technology for Human Biometric Traits Identification." *IEEE 10th International Conference on Information Technology* (2011): Published by IEEE.
  - [22] T. Ko. "Multimodal Biometric Identification for Large User Population Using Fingerprint, Face, and Iris Recognition." *IEEE Applied Imagery Pattern Recognition Workshop* (2005): Published by IEEE.
  - [23] P. P. Sarangi, D. R. Nayak, M. Panda, B. Majhi. "A Feature-Level Fusion Based Improved Multimodal Biometric Recognition System Using Ear and Profile Face." *Journal of Ambient Intelligence and Humanized Computing* (2022): Published by Springer.
  - [24] S. Chaudhary, R. Nath. "A Multimodal Biometric Recognition System Based on Fusion of Palmprint, Fingerprint, and Face." *International Conference on Advances in Recent Technologies in Communication and Computing* (2009): Published by IEEE.
  - [25] R. Shyam, Y. N. Singh. "Identifying Individuals Using Multimodal Face Recognition Techniques." *Procedia Computer Science* (2015): Published by Elsevier.
  - [26] A. Lumini, L. Nanni. "Overview of the Combination of Biometric Matchers." *Information Fusion* (2017): Published by Elsevier.
  - [27] M. He, S. J. Horng, P. Fan, R. S. Run, R. J. Chen, J. L. Lai. "Performance Evaluation of Score Level Fusion in Multimodal Biometric Systems." *Pattern Recognition* (2010): Published by Elsevier.
  - [28] M. M. H. Ali, A. T. Gaikwad. "Multimodal Biometrics Enhancement Recognition System Based on Fusion of Fingerprint and Palmprint: A Review." *Journal of Computer Science and Technology* (2016): Published by Academia.edu.
  - [29] A. A. Ross, A. K. Jain, K. Nandakumar. "Information Fusion in Biometrics." *Handbook of Multibiometrics* (2006): Published by Springer.
  - [30] A. Abdul-Al, G. K. Kyeremeh, R. Qahwaji. "A Novel Approach to Enhancing Multi-Modal Facial Recognition: Integrating Convolutional Neural Networks, Principal Component Analysis, and Sequential Neural Networks." *IEEE Access* (2024): Published by IEEE.