



Wearable Audio-Visual Enhanced Speech-recognition System (WAVESS): A Conceptual Model with Lipreading Capabilities

Research Article

<https://stem.techspherejournal.com>

Article Info

Revised Date: 21st September, 2025
Accepted Date: 23rd September, 2025
Published Date: 24th September, 2025

Keywords

Multimodal Speech Recognition
Wearable Technology
Audio-Visual Fusion
Lipreading Systems
Human-Computer Interaction

Author Details

Adeoye A.E.^{1*}, Olaye E.², Onwuegbuzie U. I.³
*1*Department of Computer Science, Dennis Osadebay University, Asaba, Delta State, Nigeria.
2 Department of Computer Engineering, University of Benin, Benin City, Edo State, Nigeria.
3 Department of Cybersecurity, Dennis Osadebay University, Asaba, Delta State, Nigeria.

*Corresponding author's email: adeoye.esther@dou.edu.ng

DOI: <https://doi.org/10.5281/zenodo.17205329>

This is an open-access article under the [CC BY-SA](https://creativecommons.org/licenses/by-sa/4.0/) license



ABSTRACT

Speech recognition technologies have evolved significantly from early rule-based systems to modern deep learning models; however, conventional audio-only approaches remain constrained by noise interference, diverse accents, and speech impairments, limiting their robustness in real-world applications. Recent research highlights the value of multimodal systems that combine auditory and visual cues, with lipreading offering complementary information where audio signals alone may fail. This study proposes the Wearable Audio-Visual Enhanced Speech-recognition System (WAVESS), a conceptual model implemented in the form of smart glasses equipped with a microphone array and a mini-camera that captures lip movements. The system integrates audio and video inputs through preprocessing pipelines, noise reduction and Mel-Frequency Cepstral Coefficients (MFCC) for audio, and lip region detection and feature extraction for video, before fusing them in a real-time multimodal recognition engine. The fused representation enhances recognition accuracy, adaptability, and resilience in challenging conditions such as noisy environments, hearing impairment contexts, and human-machine interaction scenarios. The model also incorporates connectivity features for wireless or edge-based computation and provides multimodal feedback through augmented reality overlays, audio, or haptic signals. WAVESS demonstrates the comparative advantage of wearable, multimodal systems in accessibility, communication, education, and security applications while addressing scalability and ethical considerations. The conceptual framework establishes a foundation for future prototyping, dataset expansion, and real-world deployment in advancing robust speech recognition research.

1 Introduction

1.1 Background of Speech Recognition Technologies

Speech recognition has evolved from rule-based models in the 1950s to sophisticated deep learning systems that power modern virtual assistants such as Siri, Alexa, and Google Assistant (Chiddarwar et al., 2025). Advances in machine learning, natural language processing, and computational power have enabled significant improvements in recognition accuracy and usability (Houssein et al., 2021). These systems are widely applied in accessibility, healthcare, telecommunications, security, and smart devices, creating a paradigm shift in how humans interact with machines.



Despite these advances, existing systems largely rely on audio inputs, making them vulnerable to performance degradation in certain contexts.

Although audio-only speech recognition has matured, it faces persistent challenges. Background noise, overlapping speech, and environmental disturbances often reduce recognition accuracy (Michelsanti et al., 2021). Variations in accents, dialects, and speech rates can also hinder system performance. For individuals with speech impairments or reduced vocal clarity, recognition errors become more frequent (Moore, 2021). Furthermore, in real-world scenarios such as busy streets, factories, or classrooms, conventional systems struggle to maintain reliability, highlighting the need for enhanced input modalities.

1.2 Importance of Multimodal (Audio-Visual) Approaches

Human communication is inherently multimodal; listeners rely not only on auditory signals but also on visual cues such as lip movements, facial expressions, and gestures (Krason et al., 2024). Lipreading, in particular, plays a crucial role in enhancing speech comprehension, especially in noisy environments. Integrating visual information with audio inputs provides complementary data that improves robustness and accuracy (Li et al., 2023). Recent advances in computer vision and deep learning have made real-time lipreading feasible, thereby opening opportunities for wearable devices capable of delivering multimodal speech recognition.

1.3 Research Gap and Problem Statement

While significant progress has been made in both audio and visual speech recognition, there is limited research on compact, wearable systems that seamlessly integrate the two modalities. Most existing solutions are confined to laboratory environments or rely on bulky hardware setups, making them impractical for everyday use (Diraco et al., 2023). This gap motivates the design of a wearable, real-time, audio-visual enhanced speech recognition system with lipreading capabilities, which we term **WAVESS**. The problem this study addresses is the lack of a conceptual framework for developing lightweight, user-friendly, and effective multimodal wearable devices that overcome the shortcomings of audio-only systems.

1.4 Aim and Objectives of the Study

This study aims to propose a conceptual model for a wearable audio-visual enhanced speech recognition system (WAVESS) that integrates lipreading capabilities. The specific objectives are to:

- a. Present the conceptual architecture of WAVESS, including its core components.
- b. Highlight how audio and visual data fusion can improve recognition performance in noisy and challenging environments.
- c. Explore potential applications and implications of WAVESS for accessibility, human–machine interaction, and communication technologies.

1.5 Paper Organization

The rest of the paper is structured as follows: Section 2 reviews existing literature on speech recognition, multimodal approaches, and wearable systems. Section 3 presents the conceptual framework of WAVESS, outlining its core components and data flow. Section 4 discusses methodological considerations, including data requirements, feature extraction, and learning models. Section 5 presents the WAVESS mathematical model, while Section 6 highlights potential applications across various domains. Section 7 provides a discussion of comparative advantages, limitations, and ethical considerations. Finally, Section 8 concludes the paper and suggests directions for future research.



2 Literature Review

2.1 Evolution of Speech Recognition Systems

Speech recognition has undergone a remarkable transformation over the past seven decades. Early systems of the 1950s and 1960s were rule-based, relying on phonetic pattern matching to recognize isolated words or digits (Batista, 2024). In the 1980s and 1990s, statistical models, particularly Hidden Markov Models (HMMs), became the dominant framework for continuous speech recognition due to their ability to model temporal sequences of speech (Singh et al., 2022). The 2000s saw the integration of Gaussian Mixture Models (GMMs) with HMMs, further improving accuracy (Gopinathan et al., 2024). More recently, deep learning techniques such as Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and Transformer-based architectures like DeepSpeech and wav2vec, have revolutionized the field, enabling end-to-end training and achieving near-human-level performance in controlled environments (Mienye & Swart, 2024). Despite these advances, real-world challenges like noise, accent variability, and speech impairments remain unresolved.

2.2 Audio-based Models: Strengths and Weaknesses

Audio-only speech recognition systems remain the most widely adopted due to their maturity, low hardware requirements, and relatively high accuracy under controlled conditions (Michelsanti et al., 2021). They are effective in quiet environments, have well-established datasets, and power mainstream applications such as voice assistants, transcription services, and voice-controlled IoT devices (Danka, 2024). However, their weaknesses become evident in noisy settings, where background interference leads to high error rates. Moreover, accent variations, speech disfluencies, and medical speech impairments create significant recognition challenges. These shortcomings highlight the limitations of relying solely on the auditory modality and reinforce the need for complementary cues to enhance performance.

2.3 Lipreading and Visual Speech Recognition Research

Lipreading, or visual speech recognition (VSR), involves interpreting speech from visual features of lip movements, facial expressions, and sometimes tongue or jaw motion (Battista et al., 2025). Traditional lipreading relied on handcrafted features, such as Active Appearance Models (AAMs), but deep learning has enabled more accurate and robust visual recognition (Fenghour et al., 2021). CNNs and Long Short-Term Memory (LSTM) networks have been successfully applied to lipreading tasks, with large-scale datasets like GRID, LRW (Lip Reading in the Wild), and AVSpeech advancing research (Exarchos et al., 2024). Recent transformer-based architectures have further improved performance by modeling long-term temporal dependencies. While lipreading alone is insufficient for reliable recognition due to ambiguities in visually similar phonemes like “p” vs. “b”, when combined with audio cues, it substantially enhances robustness in noisy environments (Li et al., 2023).

2.4 Wearable Technology for Human–Computer Interaction

Wearable devices such as smart glasses, hearing aids, and AR/VR headsets are redefining human–computer interaction (HCI) (Stephanidis & Salvendy, 2024). These technologies offer portability, real-time sensing, and user-centred design, enabling seamless integration of assistive and interactive functionalities into daily life. For example, Google Glass and Microsoft HoloLens have demonstrated the potential of wearable interfaces in healthcare, education, and industrial settings (Baashar et al., 2023). In speech recognition research, wearables provide an opportunity to embed microphones, cameras, and processors close to the source of speech, reducing signal distortion and latency. However, challenges related to power consumption, device ergonomics, and privacy concerns persist in deploying wearable solutions for speech recognition.



2.5 Existing Multimodal Speech Recognition Frameworks

Multimodal speech recognition integrates audio and visual signals to improve system robustness and accuracy (Jeon & Kim, 2022). Early frameworks employed feature fusion (combining audio and visual features at the input level) or decision fusion (combining outputs of separate models). Recent deep learning approaches leverage joint embeddings to model the correlation between modalities more effectively. Studies have shown that audio-visual fusion consistently outperforms unimodal systems, particularly in environments with high background noise (Jiao et al., 2024). Systems like AVSR (Audio-Visual Speech Recognition) models have achieved significant error-rate reductions compared to audio-only baselines (Ivanko et al., 2023). Nonetheless, most frameworks are implemented in controlled environments or with desktop-class computing power, limiting their real-world, wearable applicability.

2.6 Identified Gaps and Opportunities

Although the fields of speech recognition, lipreading, and wearable technology have advanced independently, their integration remains limited. Current research gaps include:

- a. A lack of lightweight, real-time wearable frameworks that fuse audio and visual modalities.
- b. Limited datasets designed for wearable, egocentric perspectives (i.e., capturing lip movements from glasses-mounted cameras).
- c. Challenges in energy efficiency, privacy, and device ergonomics that constrain adoption.
- d. Insufficient exploration of applications for accessibility, particularly for hearing-impaired individuals.

These gaps present opportunities for designing a conceptual framework like WAVESS (Wearable Audio-Visual Enhanced Speech-recognition System), which bridges speech recognition and wearable computing with lipreading capabilities. WAVESS aims to overcome the constraints of audio-only models and extend the practicality of multimodal recognition to real-world, everyday scenarios.

3 Conceptual Framework of WAVESS

3.1 Overview of Proposed Model

The Wearable Audio-Visual Enhanced Speech-recognition System (WAVESS) is conceptualized as a pair of smart glasses that seamlessly integrates audio and visual sensing for improved speech recognition. The choice of smart glasses as the wearable mode is motivated by their natural alignment with the system's functional requirements. A mini-camera embedded within the frame is oriented to capture lip movements in real time, enabling visual speech recognition (lipreading). Simultaneously, a microphone array positioned along the arms or nose bridge of the glasses captures high-fidelity audio signals, even in noisy environments. These multimodal inputs are processed by an onboard fusion engine, which may be integrated into the glasses' frame or wirelessly connected to an edge device for computational efficiency. The design further incorporates a feedback interface, allowing users to receive outputs through visual overlays on the lenses, auditory prompts, or subtle haptic signals. By leveraging the socially acceptable and ergonomic form factor of smart glasses, WAVESS ensures a practical, portable, and unobtrusive wearable solution, making it suitable for diverse real-world applications such as accessibility, industrial communication, and human-machine interaction.

3.2 System Components

a. Audio Capture Module (Microphone Array):

The microphone array captures high-quality acoustic signals, filters background noise, and extracts relevant speech features. Using beamforming and noise suppression techniques, it ensures clear signal acquisition in dynamic environments.



b. Video Capture Module (Mini-camera for Lip Movements):

A lightweight, glasses-mounted mini-camera is used to continuously track lip movements. This camera provides real-time video streams for lipreading and facial motion analysis, enabling complementary visual features to support speech recognition.

c. Processing Unit (Real-time Fusion Engine):

The processing unit is the computational hub of WAVESS. It integrates both audio and visual features, synchronizes temporal sequences, and executes machine learning models for multimodal recognition. Optimized lightweight neural networks or edge-based accelerators are employed to ensure low-latency inference.

d. Display and Feedback Interface:

WAVESS incorporates a visual or auditory feedback interface, such as a heads-up display (HUD) on smart glasses or audio feedback via bone-conduction speakers. This interface provides real-time recognition results, confirmations, or error notifications to the user.

e. Connectivity (Wireless/Edge Computing):

The system supports wireless connectivity (Wi-Fi, Bluetooth, or 5G) to offload computation when necessary or to integrate with external systems such as smartphones, cloud servers, or assistive platforms. Edge computing ensures real-time operation without over-reliance on cloud infrastructure.

3.3 Data Flow and Architecture

The WAVESS architecture follows a sequential and synchronised data flow:

- a. **Signal Acquisition:** The microphone array captures audio while the mini-camera records video of the speaker's lips.
- b. **Preprocessing:** Audio signals undergo noise reduction and spectral feature extraction using Mel-Frequency Cepstral Coefficients (MFCCs) or spectrograms, while video frames are processed to detect and extract lip regions.
- c. **Feature Extraction:** Both modalities are transformed into feature vectors. Audio features represent frequency and temporal variations, while visual features capture spatial-temporal lip motion patterns.
- d. **Fusion Engine:** Extracted features are synchronised and combined using multimodal fusion strategies.
- e. **Recognition and Output:** The fused features are processed by a recognition model to generate text or speech output, which is presented via the feedback interface.

3.4 Fusion Strategies (Feature-level vs. Decision-level)

Two principal fusion approaches can be implemented in WAVESS:

1. **Feature-level Fusion:** Audio and visual features are concatenated or combined at the input stage before being fed into a recognition model. This approach leverages joint representation learning, enabling the system to exploit inter-modal correlations. However, it requires precise synchronisation and high computational resources.
2. **Decision-level Fusion:** Independent models process audio and visual data separately, and their recognition outputs are combined through voting, weighting, or confidence-based rules. This approach offers flexibility, reduced latency, and robustness when one modality is degraded (e.g., poor lighting or high noise).

WAVESS is designed to support both strategies, with adaptive selection depending on environmental conditions and device constraints.



3.5 Conceptual Diagram of WAVESS

The conceptual framework of WAVESS can be visualised as a modular architecture consisting of:

- a. **Input Layer:** Microphone array + mini-camera.
- b. **Preprocessing Layer:** Audio filtering, lip detection, and frame preprocessing.
- c. **Feature Extraction Layer:** Acoustic and visual feature encoding.
- d. **Fusion Engine:** Multimodal integration (feature-level or decision-level).
- e. **Recognition Layer:** Multimodal deep learning model for speech decoding.
- f. **Output Layer:** Display or auditory feedback via interface, plus wireless connectivity for external integration.

Figure 1 shows the workflow of WAVESS.

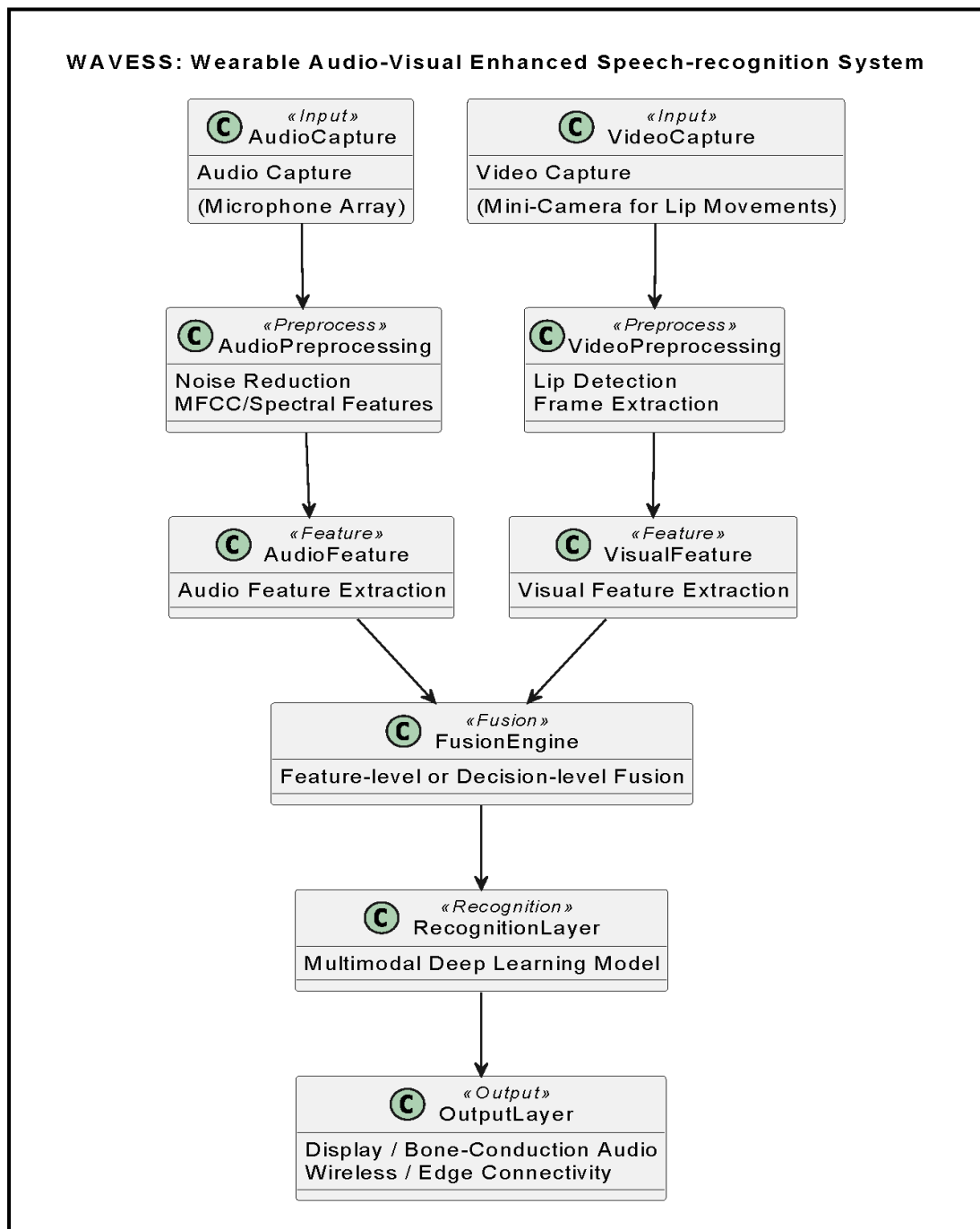


Figure 1: WAVES workflow diagram

4 Methodological Considerations

4.1 Design Approach (Wearable Device Conceptualization)

The proposed WAVESS device is conceptualized as a lightweight, glasses-mounted wearable capable of capturing both audio and visual speech signals. The design prioritizes portability, user comfort, and real-time operation. The device



integrates a microphone array for high-fidelity audio capture and a mini-camera positioned near the eyes to monitor lip movements. A compact processing unit is embedded or paired with an edge device (e.g., a smartphone) to perform real-time feature extraction, multimodal fusion, and recognition. The conceptual design also considers power efficiency, minimal latency, and unobtrusive aesthetics to encourage everyday use, especially for accessibility and assistive applications.

4.2 Data Requirements

1. Audio Datasets:

Robust audio datasets are required to train and evaluate the system under various acoustic conditions. These datasets should include:

- a. Multiple speakers with diverse accents, ages, and genders.
- b. Noisy environments to simulate real-world conditions.
- c. Both isolated words and continuous speech.

Example datasets include LibriSpeech, TED-LIUM, and VCTK.

2. Lipreading Video Datasets:

Visual data of lip movements is essential to complement audio signals. Required datasets should provide:

- a. High-resolution frontal video of speakers' lips.
- b. Variations in lighting, facial expressions, and speaking styles.
- c. Synchronized audio to support multimodal fusion.

Well-known datasets include GRID, LRW (Lip Reading in the Wild), and TCD-TIMIT.

Preprocessing may include video frame alignment, mouth region cropping, and normalisation for consistent input dimensions.

4.3 Signal Preprocessing and Feature Extraction

1. Audio Signal Preprocessing:

- a. Noise reduction using spectral subtraction or adaptive filtering.
- b. Normalization and framing of audio signals.
- c. Feature extraction such as Mel-Frequency Cepstral Coefficients (MFCCs), spectrograms, or log-Mel features.

2. Visual Signal Preprocessing:

- a. Lip detection and segmentation using computer vision techniques or deep learning-based detectors.
- b. Temporal alignment of frames with audio signals.
- c. Extraction of visual features using CNN-based embeddings or spatio-temporal modeling techniques (e.g., 3D CNNs or LSTMs).

The goal is to generate robust feature vectors from both modalities that can be fused efficiently for recognition.

4.4 Machine Learning / Deep Learning Integration

WAVESS leverages multimodal deep learning models to combine audio and visual features effectively. Possible approaches include:

- a. **Feature-level fusion:** Concatenate audio and visual feature vectors into a joint representation, followed by processing through a deep neural network (DNN), CNN-LSTM hybrid, or transformer model.
- b. **Decision-level fusion:** Independent audio and visual models generate predictions that are then combined via weighted voting or confidence-based mechanisms.



- c. **Temporal modelling:** Recurrent architectures (LSTM/GRU) or transformers capture sequential dependencies in both audio and visual streams.

The model is trained using supervised learning, with cross-entropy or CTC loss for sequence-to-sequence speech recognition tasks. Data augmentation techniques (noise injection, video transformations) are applied to improve generalization.

4.5 Performance Evaluation Criteria (Accuracy, Latency, Robustness)

Evaluation of WAVESS considers both recognition quality and practical usability:

- a. **Accuracy:** Word error rate (WER) or phoneme error rate (PER) is measured for both audio-only, video-only, and fused models.
- b. **Latency:** End-to-end processing time is monitored to ensure real-time operation suitable for wearable deployment.
- c. **Robustness:** System performance is tested under varying noise levels, lighting conditions, accents, and speaking speeds.
- d. **Energy Efficiency:** Power consumption is assessed to ensure wearable usability over extended periods.

These metrics provide a comprehensive assessment of the system's practical feasibility, reliability, and suitability for real-world applications.

Figure 2 presents the methodological process flow of the WAVESS audio-visual speech recognition pipeline.

Methodological Flow Description:

- a. **Data Acquisition:** Audio from microphone array + video from mini-camera.
- b. **Preprocessing:** Noise reduction for audio, lip detection and frame extraction for video.
- c. **Feature Extraction:** Extract audio (MFCC, spectrogram) and visual (CNN embeddings) features.
- d. **Fusion Engine:** Combine features using feature-level or decision-level fusion.
- e. **Evaluation:** Test accuracy, latency, robustness, and energy efficiency.

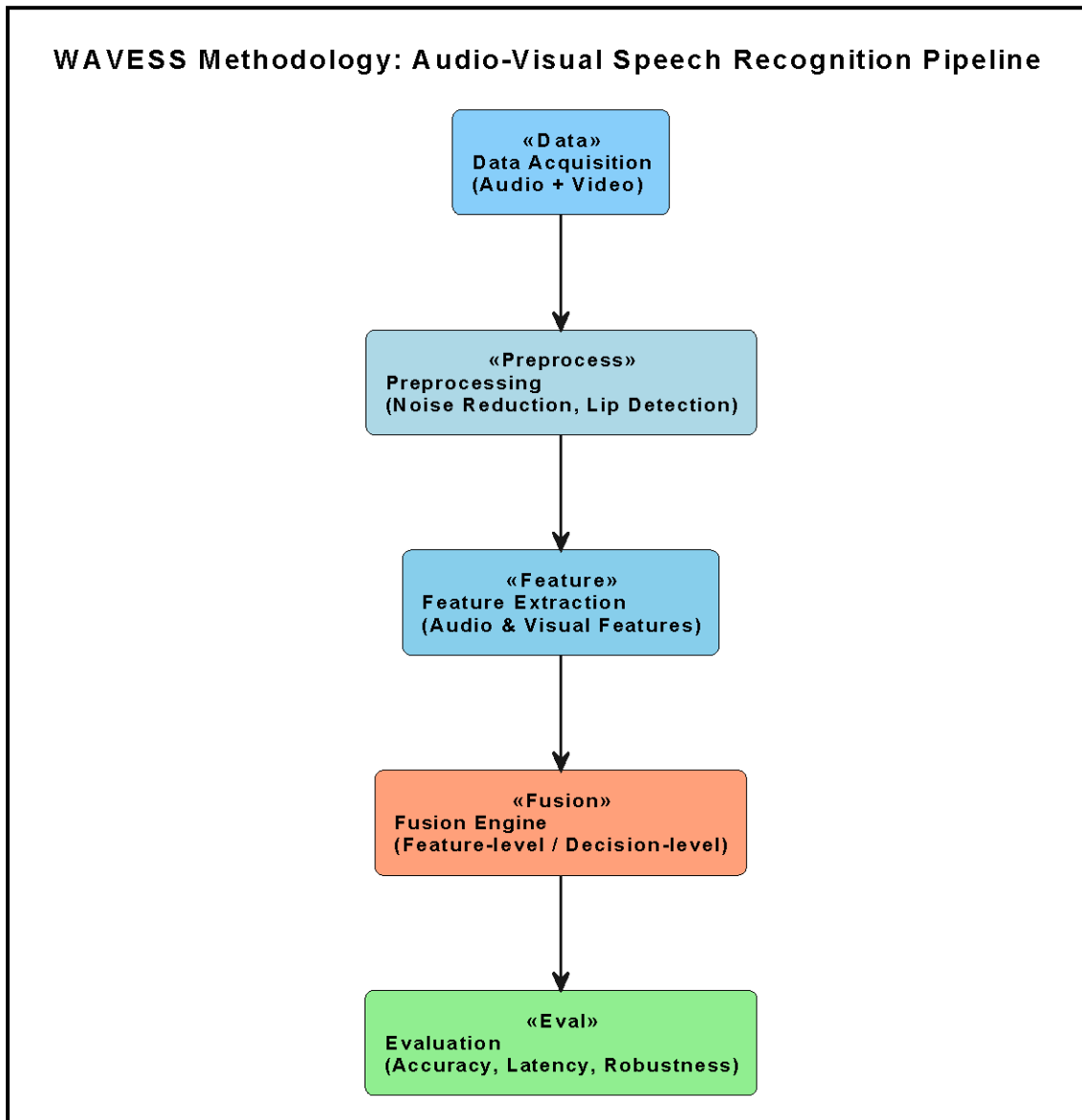


Figure 2: WAVESS audio-visual speech recognition pipeline

5 Mathematical Model of WAVESS

The WAVESS system is a multimodal audio-visual speech recognition framework, where both audio and visual (lipreading) inputs are combined to improve recognition accuracy. The system can be formally described using a mathematical model represented by Equations 1 to 7, while Equation 8 can be used to evaluate WAVESS Word Error Rate (WER) against other related models.



5.1 Audio and Visual Signal Representation

Let the audio signal captured by the microphone array be represented as:

$$A(t) = s(t) + n(t) \dots \dots \dots (1)$$

Where:

- $A(t)$ = Captured audio signal at time t
- $s(t)$ = clean speech signal
- $n(t)$ = additive noise

The visual signal (lip movements) captured by the camera is represented as a sequence of video frames:

$$V = \{v_1, v_2, v_3, \dots, v_T\} \dots \dots \dots (2)$$

Where:

- v_1 = lip region features at frame t
- T = total number of frames corresponding to the speech segment.

5.2 Feature Extraction

For audio, feature extraction (e.g., MFCC or spectrogram) transforms $A(t)$ into a feature vector:

$$f_a = F_a(A) \in \mathbb{R}^{d_a} \dots \dots \dots (3)$$

Where:

- F_a = audio feature extraction function
- d_a = dimensionality of audio features

For visual input, lip region frames are processed to obtain visual features:

$$f_v = F_v(V) \in \mathbb{R}^{d_v} \dots \dots \dots (4)$$

Where:

- F_v = visual feature extraction function (e.g., CNN embedding)
- d_v = dimensionality of visual features

5.3 Multimodal Fusion

The fusion module integrates audio and visual features to form a joint representation:

$$f_{av} = \Phi(f_a, f_v) \dots \dots \dots (5)$$

Where:

- f_{av} = combined audio-visual feature vector
- Φ = fusion function, which can be:
 - **Feature-level fusion:** $f_{av} = [f_a; f_v]$ (concatenation)



- **Decision-level fusion:** $y = w_a y_a + w_v y_v$, where y_a, y_v are predictions from audio and visual models, and w_a, w_v are confidence weights.

5.4 Recognition Function

The fused feature vector is input to a recognition model R (e.g., deep neural network, LSTM, or transformer) to produce the predicted transcription:

$$\hat{y} = R(f_{av}) \dots \dots \dots (6)$$

Where:

\hat{y} = predicted text sequence

R = mapping function learned via supervised training

The objective during training is to minimize the loss function, commonly cross-entropy or Connectionist Temporal Classification (CTC) loss:

$$L = Loss(\hat{y}, y_{true}) \dots \dots \dots (7)$$

Where y_{true} is the ground-truth transcription.

5.5 System Performance Metrics

The effectiveness of WAVESS can be evaluated using:

1. **Word Error Rate (WER):**

$$WER = \frac{S+D+I}{N} \times 100\% \dots \dots \dots (8)$$

Where S = substitutions, D = deletions, I = insertions, and N = total words.

2. **Latency (L):** Time delay from input acquisition to output generation.
3. **Robustness (R_b):** Accuracy under varying noise levels and visual conditions.

5.6 Explanation of the Model

- 1 **Audio Module:** Captures speech signals and extracts relevant acoustic features.
- 2 **Visual Module:** Captures lip movements and extracts visual features.
- 3 **Fusion Module:** Combines audio and visual features either at the feature or decision level, improving recognition robustness in noisy or challenging environments.
- 4 **Recognition Module:** Maps fused features to textual output using deep learning models.
- 5 **Performance Evaluation:** Quantifies system reliability, accuracy, and real-time applicability.

This mathematical formulation provides a formal representation of WAVESS, highlighting how audio-visual fusion can systematically enhance speech recognition performance. Figure x presents a compact diagrammatic representation of the mathematical flow, linking **audio** → **visual** → **fusion** → **recognition** → **evaluation**

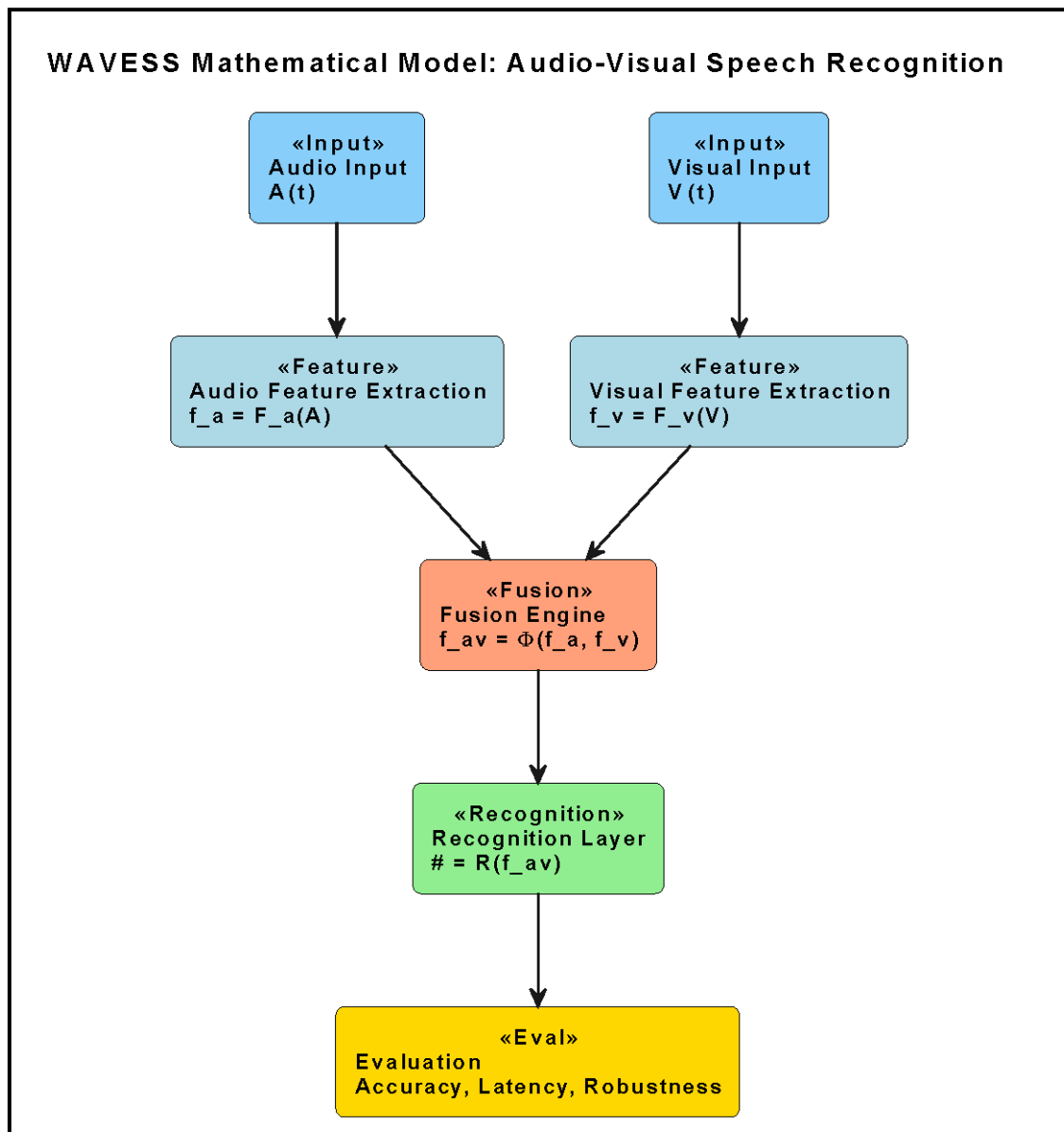


Figure 3: WAVESS mathematical model diagram

6 Potential Applications

6.1 Accessibility for Hearing-Impaired Users

One of the most significant applications of WAVESS is enhancing communication for individuals with hearing impairments. By combining audio and visual speech recognition, the system can provide real-time transcription of spoken words, displayed through a heads-up display or audio feedback devices. Lipreading capabilities enable users to understand speech even in cases where audio signals are weak or distorted. This functionality can be integrated into daily life scenarios such as meetings, lectures, and social interactions, greatly improving accessibility and social inclusion for the hearing-impaired population.



6.2 Enhanced Communication in Noisy Environments

Traditional audio-only speech recognition systems struggle in environments with high background noise, such as factories, construction sites, traffic-heavy areas, or crowded public spaces (Jeon & Kim, 2022). WAVESS addresses this challenge by leveraging visual cues from lip movements, which are unaffected by ambient noise. The multimodal approach ensures reliable speech recognition in acoustically challenging conditions, allowing professionals to communicate effectively, maintain operational efficiency, and reduce the likelihood of errors in critical tasks.

6.3 Human–Machine Interaction (Assistive Devices, Smart Assistants)

WAVESS can significantly enhance human–machine interaction by providing more robust voice-controlled interfaces. The integration of lipreading improves command recognition for smart assistants, IoT devices, and wearable AI assistants, even in noisy surroundings or when the user speaks softly. Furthermore, the wearable design allows seamless interaction in real-world scenarios without relying on desktop systems, supporting applications in personal assistive devices, robotics control, and augmented reality environments.

6.4 Security and Authentication

Audio-visual speech recognition can strengthen security and authentication systems. WAVESS can be employed in biometric authentication by verifying both vocal and lip movement patterns, creating a multimodal biometric signature that is difficult to spoof. This approach can be applied to secure device access, financial transactions, or identity verification in sensitive environments. The dual-modality reduces false positives and enhances overall system reliability compared to unimodal audio-based authentication.

6.5 Education and Language Learning

In educational settings, WAVESS can assist both teachers and learners by providing real-time speech-to-text transcription, pronunciation guidance, and visual feedback on lip movements. Language learners can benefit from lipreading analysis to improve pronunciation and comprehension. Additionally, the system can support classroom accessibility for hearing-impaired students, enabling them to follow lectures and participate fully in discussions. These capabilities make WAVESS a valuable tool for inclusive, technology-enhanced learning environments.

7 Discussion

7.1 Comparative Advantage of WAVESS over Traditional Systems

The primary advantage of WAVESS lies in its multimodal design, which integrates both audio and visual information. Traditional audio-only speech recognition systems suffer significant performance degradation in noisy environments or in cases of speech impairments. By incorporating lipreading capabilities, WAVESS provides greater robustness, accuracy, and reliability, even in acoustically challenging scenarios such as crowded public spaces or industrial sites. Moreover, its wearable implementation enhances portability and accessibility, making it suitable for real-time, on-the-go communication. This dual-modality approach represents a substantial improvement over conventional systems, aligning with current trends toward more inclusive and resilient human–computer interaction technologies.



7.2 Expected Limitations

Despite its promising features, WAVESS is not without challenges. The use of a mini-camera for lipreading raises privacy concerns, as continuous video capture may inadvertently record sensitive environments or individuals. Power consumption is another limitation, as simultaneous audio-visual processing, real-time fusion, and wireless connectivity demand significant energy, potentially affecting battery life in wearable devices. Furthermore, the system's accuracy may decline in conditions with poor lighting or occluded faces, limiting its performance in certain real-world settings. These limitations highlight the need for optimization in both hardware and algorithms to achieve practical, long-term usability.

7.3 Possible Improvements and Scalability

To address the limitations, several improvements can be envisioned. Integration of low-power edge AI processors could significantly reduce latency and energy demands. Advanced privacy-preserving techniques, such as on-device processing and encrypted data handling, would help mitigate security concerns. Scalability could be achieved by designing WAVESS to support multiple languages, accents, and dialects, thereby broadening its applicability across diverse populations. Furthermore, incorporating adaptive learning mechanisms could enable the system to personalize recognition performance for individual users over time. These enhancements would ensure that WAVESS can evolve from a conceptual framework to a scalable solution suitable for mass deployment.

7.4 Ethical and Social Implications

The ethical and social dimensions of WAVESS are critical considerations. On one hand, the system offers inclusive technologies that empower hearing-impaired individuals, enhance workplace safety, and improve educational access. On the other hand, issues of data privacy, surveillance risks, and potential misuse must be carefully addressed. There is also a need to ensure fair accessibility and affordability, so that such innovations do not widen the digital divide. Transparent guidelines, regulatory compliance, and responsible AI practices will be essential to maximize the positive social impact of WAVESS while minimizing risks.

8 Conclusion and Future Work

8.1 Summary of Key Contributions

This paper presents WAVESS (Wearable Audio-Visual Enhanced Speech Recognition System) as a conceptual model that integrates audio and lip-reading-based visual cues to improve the accuracy and robustness of speech recognition. The study provided a theoretical framework, mathematical formulation, and methodological design that highlight how multimodal fusion can overcome the limitations of conventional audio-only systems, particularly in noisy environments and for individuals with speech or hearing challenges. By outlining the system's architecture, potential applications, and ethical considerations, the paper contributes a comprehensive roadmap for advancing wearable multimodal speech recognition technologies.

8.2 Impact of WAVESS in Advancing Speech Recognition Research

WAVESS underscores the transformative role of multimodal approaches in bridging the performance gaps of current speech recognition systems. Its integration of wearable technology, audio-visual fusion, and real-time processing offers a new paradigm for human-computer interaction. Beyond academic contribution, WAVESS has the potential to influence practical domains such as accessibility, security, education, and industrial communication, thereby reinforcing the societal and technological significance of this research direction.



8.3 Future Research Directions

- a. **Prototype Development:** Future efforts should focus on translating the conceptual model into a functional prototype. This will involve hardware integration of microphones, mini-cameras, and processing units in a lightweight wearable design, alongside software implementation of the multimodal recognition framework.
- b. **Dataset Expansion:** Current datasets for multimodal speech recognition are often limited in terms of language diversity, accents, and environmental conditions. Expanding datasets to include multilingual, multi-accent, and noisy real-world recordings will be essential for training more robust models.
- c. **Real-World Trials:** Extensive field evaluations will be necessary to validate WAVESS in diverse contexts such as classrooms, factories, outdoor spaces, and healthcare settings. These trials will help assess usability, reliability, and user acceptance, providing valuable feedback for iterative improvement of both hardware and software.

References

- Baashar, Y., Alkaws, G., Wan Ahmad, W. N., Alomari, M. A., Alhussian, H., & Tiong, S. K. (2023). Towards wearable augmented reality in healthcare: a comparative survey and analysis of head-mounted displays. *International Journal of Environmental Research and Public Health*, 20(5), 3940.
- Batista, J. R. (2024). *Learn OpenAI Whisper: Transform your understanding of GenAI through robust and accurate speech processing solutions*. Packt Publishing Ltd.
- Battista, M., Collese, F., Orzan, E., Fantoni, M., & Bottari, D. (2025). Lip-Reading: Advances and Unresolved Questions in a Key Communication Skill. *Audiology Research*, 15(4), 89.
- Chidharwar, G., Bhabad, S., & Indalkar, A. (2025). Advancements in Conversational AI: A Survey of Chatbots and Voice Assistants. *2025 1st International Conference on AIML-Applications for Engineering & Technology (ICAET)*, 1–8.
- Danka, S. (2024). *The effectiveness of a Computer-Assisted Pronunciation Training Approach for the production and perception of linking by English L2 learners*. University of Leicester.
- Diraco, G., Rescio, G., Siciliano, P., & Leone, A. (2023). Review on human action recognition in smart living: Sensing technology, multimodality, real-time processing, interoperability, and resource-constrained processing. *Sensors*, 23(11), 5281.
- Exarchos, T., Dimitrakopoulos, G. N., Vrahatis, A. G., Chrysovitiotis, G., Zachou, Z., & Kyrodimos, E. (2024). Lip-reading advancements: A 3D convolutional neural network/long short-term memory fusion for precise word recognition. *BioMedInformatics*, 4(1), 410–422.
- Fenghour, S., Chen, D., Guo, K., Li, B., & Xiao, P. (2021). Deep learning-based automated lip-reading: A survey. *IEEE Access*, 9, 121184–121205.
- Gopinathan, K. N., Murugesan, P., & Jeyaraj, J. J. (2024). Stock price prediction using a novel approach in Gaussian mixture model-hidden Markov model. *International Journal of Intelligent Computing and Cybernetics*, 17(1), 61–100.
- Houssein, E. H., Mohamed, R. E., & Ali, A. A. (2021). Machine learning techniques for biomedical natural language processing: a comprehensive review. *IEEE Access*, 9, 140628–140653.
- Ivanko, D., Ryumin, D., & Karpov, A. (2023). A review of recent advances on deep learning methods for audio-visual speech recognition. *Mathematics*, 11(12), 2665.
- Jeon, S., & Kim, M. S. (2022). Noise-robust multimodal audio-visual speech recognition system for speech-based interaction applications. *Sensors*, 22(20), 7738.
- Jiao, T., Guo, C., Feng, X., Chen, Y., & Song, J. (2024). A Comprehensive Survey on Deep Learning Multi-Modal Fusion: Methods, Technologies and Applications. *Computers, Materials & Continua*, 80(1).
- Krason, A., Varley, R., & Vigliocco, G. (2024). Understanding discourse in face-to-face settings: The impact of multimodal cues and listening conditions. *Journal of Experimental Psychology: Learning, Memory, and Cognition*.
- Li, D., Gao, Y., Zhu, C., Wang, Q., & Wang, R. (2023). Improving speech recognition performance in noisy environments by enhancing lip reading accuracy. *Sensors*, 23(4), 2053.
- Michelsanti, D., Tan, Z.-H., Zhang, S.-X., Xu, Y., Yu, M., Yu, D., & Jensen, J. (2021). An overview of deep-learning-based audio-visual speech enhancement and separation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29, 1368–1396.
- Mienye, I. D., & Swart, T. G. (2024). A comprehensive review of deep learning: Architectures, recent advances, and applications. *Information*, 15(12), 755.
- Moore, M. (2021). Speech Recognition for Individuals with Voice Disorders. In *Multimedia for Accessible Human Computer Interfaces* (pp. 115–144). Springer.
- Singh, A., Kaur, N., Kukreja, V., Kadyan, V., & Kumar, M. (2022). Computational intelligence in processing of speech acoustics: a survey. *Complex & Intelligent Systems*, 8(3), 2623–2661.
- Stephanidis, C., & Salvendy, G. (2024). *Interaction Techniques and Technologies in Human-Computer Interaction*. CRC Press.