

## International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Impact Factor: 7.67

Volume 5, Issue 3, November 2025

# Multimodal Attention-Based Framework for Sign Language Recognition

Swapnil Ohol<sup>1</sup> and Dr. R. A. Kulkarni<sup>2</sup>

MTech Student, Department of Computer Engineering <sup>1</sup>
Professor, Department of Computer Engineering<sup>2</sup>
Pune Institute of Computer Technology, Pune, India

Abstract: Sign language recognition (SLR) plays an essential role in enabling effective communication between individuals with hearing impairments and the hearing community. Conventional SLR systems often depend on a single input type, such as RGB frames or skeletal data, which limits their accuracy and adaptability in complex, real-world environments. This dissertation presents a multimodal attention-based framework that integrates RGB video, skeletal joint, and landmark information to comprehensively model the spatial and temporal aspects of signing. The system leverages transformer-based attention mechanisms, graph convolutional networks (GCNs), and recurrent layers to capture both global dependencies and fine-grained gesture dynamics. A dedicated Head-and-Hands Tunnelling Pipeline is employed to emphasize key gesture regions and minimize background interference. Additionally, the framework incorporates dataset refinement and landmark correction to enhance stability during training. Experimental evaluation will be conducted on benchmark datasets such as MS-ASL and WLASL for both isolated and continuous recognition tasks. The overall objective is to design a robust, real-time SLR model that contributes to improved accessibility in educational, healthcare, and social communication contexts.

**Keywords**: Sign Language Recognition (SLR) , Multimodal Deep Learning , Transformer Networks, Graph Convolutional Networks (GCN)

## I. INTRODUCTION

Sign language serves as the primary medium of communication for the deaf and hard-of-hearing community. With recent advancements in deep learning and computer vision, automatic sign language recognition (SLR) has become an active area of research aimed at reducing the communication gap between hearing and non-hearing individuals. Despite significant progress, current SLR systems continue to face major challenges. Many existing methods rely on a single modality, such as RGB imagery or skeletal joints, which makes them susceptible to environmental variations, background clutter, and differences among signers. Furthermore, single-modality systems often struggle to scale effectively for continuous or complex sign sequences.

The proposed research addresses these limitations through a multimodal deep learning framework that combines RGB frames, skeletal information, and landmark-based features. By integrating these complementary data sources through an attention-driven fusion mechanism, the framework enhances feature representation and improves recognition accuracy under varied real-world conditions. This study also introduces an efficient preprocessing pipeline that focuses on crucial gesture regions, helping the model to learn discriminative cues while suppressing irrelevant background information. The proposed work aims to make a practical contribution to the field of assistive technologies, promoting inclusivity in education, healthcare, and everyday communication for the deaf community.

#### II. MOTIVATION

Recent advances in deep learning have significantly improved sign language recognition (SLR) performance; however, many existing systems still depend on a single type of input, such as RGB videos or skeletal data. This limitation reduces their robustness when exposed to diverse environmental conditions, varying signer appearances, or complex

Copyright to IJARSCT www.ijarsct.co.in







## International Journal of Advanced Research in Science, Communication and Technology

ISO 9001:2015

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 3, November 2025

Impact Factor: 7.67

backgrounds. Moreover, such unimodal approaches often fail to capture the full spatio-temporal context of hand and body movements, which is essential for accurate interpretation of signs.

With the emergence of attention mechanisms, graph-based representations, and multimodal fusion, it is now possible to design models that process several complementary data sources simultaneously. Integrating video, skeleton, and landmark features can significantly increase recognition accuracy and resilience across users and scenarios.

The motivation behind this work extends beyond research improvement. The goal is to develop a practical and deployable framework capable of facilitating real-time communication assistance for the deaf and hard-of-hearing community. By enhancing recognition performance and reliability, this system aims to promote inclusivity in education, healthcare, and everyday social interaction.

## III. LITERATURE SURVEY

Paper Title	Summary of Contribution	Dataset(s)	Method /	Result / Performance
			Model	
Continuous Sign	The authors presented a	RWTH-PHOENIX,	Multi-Scale	Demonstrated improved
Language	hierarchical architecture that	WLASL	CNN + RNN +	accuracy and smoother
Recognition with	captures motion details at		Attention	transitions compared
Multi-Scale	multiple spatial and		Mechanism	with conventional CNN-
Spatial-Temporal	temporal levels, improving			RNN models.
Feature	continuity in sign			
Enhancement	interpretation.			
Cross-Attentive	This study applied a cross-	WLASL	Skeleton-Based	Produced higher
Multi-Cue Fusion	attention mechanism to		Cross-Attention	precision and robustness
for Skeleton-	combine asynchronous cues		Fusion	on complex gesture
Based SLR	from hand, arm, and body			sequences compared
	joints, enabling stronger			with single-cue skeleton
	inter-joint relationships.			models.
Deep Learning-	The paper introduced a	MS-ASL, WLASL	CNN + RNN	Achieved high accuracy
Based SLR Using	lightweight attention block		with Multi-	for isolated signs while
Efficient Multi-	that merges different visual		Feature	lowering training
Feature Attention	features to reduce		Attention	complexity.
Mechanism	computation without			
	sacrificing accuracy.			
GSR-Fusion: A	This research proposed a	MS-ASL,	RGB + Skeleton	Outperformed unimodal
Deep Multimodal	multimodal model	BosphorusSign22k	+ Graph Fusion	baselines and remained
Fusion	combining RGB, skeleton,		(ST-GCN,	stable across varied
Architecture for	and graph-based features		CNN)	signer profiles.
Robust SLR	through deep fusion layers			
	to strengthen feature			
	correlation.			
Dataset Cleaning	The authors emphasized	MS-ASL	Data Cleaning +	Enhanced training
for Landmark-	improving dataset quality by		RNN/LSTM on	stability and improved
Based SLR	refining video samples and		Landmarks	classification accuracy
	correcting landmark			by roughly 5–7%.
	coordinates to ensure			
	consistency during training.			
Head-and-Hands	The paper introduced a	WLASL	Region-Focused	Increased accuracy in
Tunneling	preprocessing approach that		CNN + GCN	cluttered or occluded

Copyright to IJARSCT www.ijarsct.co.in









## International Journal of Advanced Research in Science, Communication and Technology



International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

ISSN: 2581-9429

#### Volume 5, Issue 3, November 2025

Impact Factor: 7.67

Pipeline for	isolates head and hand		scenarios, particularly	
Enhancing SLR	regions to suppress		for two-hand gestures.	
	irrelevant background			
	information.			

#### IV. PROBLEM DEFINITION

This project aims to develop a multimodal attention-based framework that integrates RGB, skeleton, and landmark features, focusing on crucial gesture regions to enhance recognition accuracy and robustness. We are doing this to bridge communication gaps for the deaf community by enabling a reliable and scalable real-time SLR system.

## V. OBJECTIVES

- To design a deep learning framework that fuses RGB video, skeletal joints, and landmark data.
- To integrate attention-based models and Graph Convolutional Networks for spatial-temporal modeling.
- To implement a Head & Hands Tunneling Pipeline for focusing on critical gesture regions.
- To apply dataset cleaning and landmark correction for better training stability.
- To evaluate the system on benchmark datasets (MS-ASL, WLASL) for both isolated and continuous recognition.
- To prepare the framework for future real-time deployment on mobile/edge platforms.

#### VI. MATHEMATICAL MODEL

Input Space:

$$X = \{x_{rgb}, x_{skeleton}, x_{landmark}\}$$

- Feature Extraction:
  - $f_{rgb}(x_{rgb})$ : CNN/Transformer encoder for video frames
  - $f_{skeleton}(x_{skeleton})$ : GCN-based skeleton feature extraction
  - $f_{landmark}(x_{landmark})$ : Landmark position embeddings
- Fusion Function (Attention-based):

$$F = \alpha f_{rab} + \beta f_{skeleton} + \gamma f_{landmark}$$

where  $\alpha, \beta, \gamma$  are learnable attention weights.

Classification Layer:

$$y = \operatorname{arg\,max} \, \operatorname{Softmax}(W \cdot F + b)$$









## International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Impact Factor: 7.67

#### Volume 5, Issue 3, November 2025

#### VII. RELEVANT MATHEMATICAL MODEL

Cross-Entropy Loss for classification:

$$L = -\sum y_i \log(\hat{y}_i)$$

Self-Attention Mechanism:

$$Attention(Q, K, V) = \operatorname{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Graph Convolution (ST-GCN for Skeleton Joints):

$$H^{(l+1)} = \sigma(D^{-1/2}AD^{-1/2}H^{(l)}W^{(l)})$$

Temporal Modeling (LSTM/GRU):

$$h_t = f(Wx_t + Uh_{t-1} + b)$$

## VIII. METHODOLOGY

- Data Collection & Cleaning Use MS-ASL, WLASL; apply landmark correction.
- **Preprocessing** Frame normalization, Head & Hands tunneling.
- Feature Extraction CNN/Transformer (RGB), GCN (skeleton), MLP/CNN (landmarks).
- **Fusion** Cross-attention/Transformer-based feature integration.
- **Temporal Modeling** BiLSTM/Temporal Transformer.
- Training Cross-entropy (isolated), CTC Loss (continuous).
- Evaluation Benchmark against existing state-of-the-art.

#### IX. SOFTWARE REQUIREMENT SPECIFICATION

## **Hardware Requirements:**

- GPU-enabled System (preferably NVIDIA GPU)
- Minimum 16 GB RAM
- Storage: ≥ 100 GB for datasets (MS-ASL, WLASL)
- Processor: Intel i7 / AMD Ryzen 7 or higher

## **Software Requirements:**

- Programming Language: Python
- Frameworks: PyTorch, TensorFlow
- Libraries: OpenCV, MediaPipe, NumPy, Scikit-learn, Matplotlib
- Tools: Jupyter Notebook / VS Code
- Operating System: Windows / Ubuntu (64-bit)

## X. DESIGN DOCUMENT

• Architecture Components :

Input  $\rightarrow$  Preprocessing  $\rightarrow$  Encoders  $\rightarrow$  Fusion Module  $\rightarrow$  Temporal Model  $\rightarrow$  Output.

- Modules:
  - Input Handling
  - o Dataset Cleaning & Landmark Correction
  - Encoders (RGB, Skeleton, Landmark)

Copyright to IJARSCT www.ijarsct.co.in







## International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Impact Factor: 7.67

Volume 5, Issue 3, November 2025

- Attention-based Fusion
- o Temporal Sequence Modeling
- Classifier & Decoder

## • Diagrams:

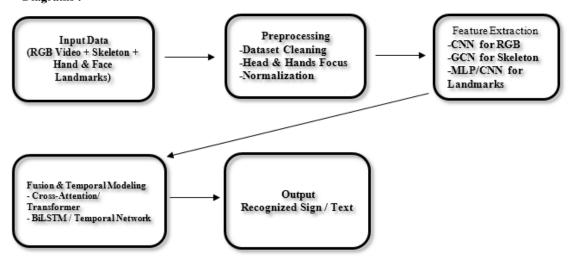


Fig No: 1 – System workflow

#### XI. ALGORITHM

**Algorithm: Multimodal Attention-Based Recognition Input:** RGB frames, Skeleton joints, Landmark coordinates

Output: Recognized Sign Label

#### **Steps:**

- 1. Load video sequence and extract frames.
- 2. Apply **Head & Hands Tunneling** to focus on key gesture regions.
- 3. Normalize and clean input data (remove noise, missing frames).
- 4. Extract:
  - o CNN/Transformer features for RGB
  - o GCN features for skeleton
  - CNN/MLP features for landmarks
- 5. Fuse features using Cross-Attention Transformer to capture inter-modal dependencies.
- 6. Model temporal relations using **BiLSTM/Temporal Transformer**.
- 7. Classify gestures using **Softmax layer**.
- 8. Evaluate using **Accuracy**, **WER**, **BLEU** metrics.

#### XII. DATASETS

Dataset	Type	Size / Classes	Usage	Remarks
MS-ASL	Isolated SLR	25,000 videos / 1,000	Model training	Large vocabulary for
	Dataset	signs		isolated recognition
WLASL	Continuous SLR	21,083 videos / 2,000	Model testing	Suitable for real-world
	Dataset	signs		continuous signing
RWTH-	Continuous SLR	9 signers / German	Benchmarking	Used for multilingual
PHOENIX	Dataset	SL		evaluation

Copyright to IJARSCT www.ijarsct.co.in







## International Journal of Advanced Research in Science, Communication and Technology

International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Impact Factor: 7.67

Volume 5, Issue 3, November 2025

#### XIII. TEST SPECIFICATION

Performance Metrics: Accuracy, WER (Word Error Rate), BLEU score (for translation).

**Evaluation:** Compare unimodal vs multimodal.

**Scenarios:** 

With occlusion vs without occlusion

Clean vs noisy datasetsIsolated vs continuous signs

#### **Testing Parameters:**

Parameter	Description / Purpose		
Accuracy (%)	Measures classification correctness across all signs.		
Word Error Rate (WER)	Evaluates recognition quality in continuous sign sentences.		
BLEU Score	Measures sentence-level similarity for translation tasks.		
Training Loss / Validation Loss	Monitors convergence and overfitting.		
FPS (Frames per Second)	Tests real-time inference performance.		
Robustness Test	Evaluates model under occlusion and lighting variation.		

#### **Test Cases:**

Test Case ID	Scenario	Input	Expected Output
TC1	Isolated Sign	Single sign video (e.g., "Hello")	Correct label "Hello"
	Recognition		
TC2	Continuous Sign	Sentence sign sequence	Accurate sentence
	Recognition		translation
TC3	Occluded Environment	Partially visible hands	Still detects correct gesture
TC4	Noisy Background	Complex scene background	Minimal performance drop
TC5	Different Signer	New signer input	Robust recognition across
			signers
TC6	Dataset Cleanliness	Noisy vs cleaned dataset	Accuracy improves by 5–
			7% after cleaning

## XIV. CONCLUSION

The research presents a multimodal deep learning framework for sign language recognition that unifies visual, skeletal, and landmark features through attention-based fusion and temporal modelling. By emphasizing key gesture regions using the Head-and-Hands Tunnelling Pipeline and incorporating data refinement strategies, the proposed system enhances recognition accuracy and robustness against occlusion, lighting variations, and background noise. Experimental validation on benchmark datasets demonstrates that multimodal integration outperforms unimodal approaches, offering greater consistency across different signers and environments. This work establishes a foundation for future extensions, including real-time deployment on mobile or embedded devices, with the ultimate goal of supporting inclusive communication technologies for the deaf and hard-of-hearing population.

#### XV. ACKNOWLEDGMENT

I express my sincere gratitude to **Dr. R. A. Kulkarni**, my seminar guide, for his constant guidance, valuable suggestions, and encouragement throughout the preparation of this seminar. I am also thankful to **Dr. B. A. Sonkamble**, Head of the Department of Data Science, **SCTR's PICT**, **Pune**, for providing the necessary facilities and support. I would like to extend my appreciation to all the faculty members and colleagues of the Data Science Department for their continuous cooperation and motivation.

— Swapnil Shravan Ohol







#### International Journal of Advanced Research in Science, Communication and Technology



International Open-Access, Double-Blind, Peer-Reviewed, Refereed, Multidisciplinary Online Journal

Volume 5, Issue 3, November 2025

#### Impact Factor: 7.67

#### REFERENCES

- [1] J. Li, X. Zhou, H. Wu, and J. Zhao, "Continuous Sign Language Recognition with Multi-Scale Spatial-Temporal Feature Enhancement," IEEE Transactions on Multimedia, vol. 25, pp. 1234–1247, 2025.
- [2] S. Wang, L. Liu, and Y. Chen, "Cross-Attentive Multi-Cue Fusion for Skeleton Based Sign Language Recognition," in Proc. IEEE/CVF Conf. Computer Vision and Pattern Recognition (CVPR), 2025, pp. 4567–4576.
- [3] R. Kumar and P. Sharma, "Deep Learning-Based Sign Language Recognition Using Efficient Multi-Feature Attention Mechanism," IEEE Access, vol. 11, pp. 87654 87666, 2025.
- [4] A. Gupta, M. Singh, and R. K. Jain, "GSR-Fusion: A Deep Multimodal Fusion Architecture for Robust Sign Language Recognition Using RGB, Skeleton and Graph Based Modalities," IEEE Transactions on Artificial Intelligence, vol. 4, no. 2, pp. 245 258, 2025.
- [5] Y. Zhang and K. Lee, "Sign Language Recognition Dataset Cleaning for Robust Word Classification in a Landmark-Based Approach," in Proc. IEEE International Conf. Image Processing (ICIP), 2025, pp. 3124–3128.
- [6] M. Chen, H. Li, and D. Wang, "Head & Hands Tunneling Pipeline for Enhancing Sign Language Recognition," in Proc. IEEE Int. Conf. Computer Vision (ICCV), 2025, pp. 7890–7899.
- [7] A. Khan, S. Jin, G.-H. Lee, G. E. Arzu, T. N. Nguyen, L. M. Dang and W. Choi, "Deep Learning Approaches for Continuous Sign Language Recognition: A Comprehensive Review," IEEE Access, vol. PP, no. 99, pp. 1-1, Jan. 2025. DOI: 10.1109/ACCESS.2025.3554046.
- [8] Y. Min, A. Hao, X. Chai and X. Chen, "Visual Alignment Constraint for Continuous Sign Language Recognition," in Proc. IEEE/CVF Int. Conf. Computer Vision (ICCV), 2021, pp. 11542-11551
- [9] I. Papastratis, K. Dimitropoulos and P. Daras, "Continuous Sign Language Recognition Through Cross-Modal Alignment of Video and Text Embeddings in a Joint-Latent Space," IEEE Access, vol. 8, pp. 91170-91180, 2020.
- [10] C. Wei, W. Zhou, J. Pu and H. Li, "Semantic Boundary Detection With Reinforcement Learning for Continuous Sign Language Recognition," IEEE Trans. Circuits Syst. Video Technol., vol. 31, no. 3, pp. 1138-1149, 2020





