

# An Integrated CNN–Hand Landmark Feature Framework for Multi-Language Sign Language Recognition

**Dr. S. Karimulla Basha**  
Dept. of CSE (DATA SCIENCE)  
RGM CET  
Nandyal, India  
kareem768@gmail.com

**Bellamkonda Masthan Basha**  
Dept. of CSE (DATA SCIENCE)  
RGM CET  
Nandyal, India  
bellamkondamasthanbasha@gmail.com

**Vadla Meghana**  
Dept. of CSE (DATA SCIENCE)  
RGM CET  
Nandyal, India  
vadlameghanavm@gmail.com

**Yangala Shashidhar Reddy**  
Dept. of CSE (DATA SCIENCE)  
RGM CET  
Nandyal, India  
yangalashashidharreddy@gmail.com

**Abstract**—The use of sign language recognition systems is critical in enhancing communication among the deaf and hard-of-hearing. This project expands an existing multi-modal American Sign Language (ASL) recognition framework into a multi-language scalable sign language recognizer. The proposed method integrates Convolutional Neural Network (CNN)-based visual feature extraction with hand landmark-based structural features to improve the accuracy and robustness of gesture classification. By combining image representations and skeletal features, the system effectively handles challenges such as variations in lighting, background noise, and hand orientation.

Unlike traditional systems that are limited to a single language, the proposed framework supports multiple sign languages, including American Sign Language (ASL), Indian Sign Language (ISL), and Arabic Sign Language (ArSL). A modular architecture is designed to ensure flexibility and scalability, where separate trained models are maintained for each language. A user-friendly interface enables dynamic language selection, allowing the system to load and execute the appropriate model in real time.

The system processes live webcam input and translates recognized gestures into text output, enabling seamless communication. This design not only enhances accuracy but also improves usability and adaptability in real-world scenarios. Overall, the proposed system presents an efficient and practical solution for multilingual sign language recognition, making it suitable for diverse cultural and communication environments.

**Index Terms**—Sign Language Recognition, Convolutional Neural Networks (CNN), Hand Landmark Detection, Multi-Modal Learning, Real-Time Gesture Recognition, American Sign Language (ASL), Indian Sign Language (ISL), Arabic Sign Language (ArSL).

## I. INTRODUCTION

Sign language is a vital mode of communication among people with hearing and speech impairments, enabling them to convey information and interact with society. However, a significant communication gap still exists between sign

language users and non-users, creating the need for automated sign language recognition systems. In recent years, deep learning techniques, particularly Convolutional Neural Networks (CNNs), have gained significant attention in this domain due to their ability to effectively extract spatial features from images [1], [2]. These models have shown promising performance in recognizing hand gestures in controlled environments and have become the foundation for many modern recognition systems [3].

Despite their effectiveness, CNN-based approaches are sensitive to variations in lighting conditions, background noise, and hand orientation, which can reduce performance in real-world scenarios. To address these limitations, researchers have explored sequential models such as Long Short-Term Memory (LSTM) networks that capture temporal dependencies in continuous gesture sequences [4]. While CNN-LSTM models improve dynamic gesture recognition, they introduce higher computational complexity and may not be suitable for real-time applications focused on single gesture classification.

More recently, hand landmark-based approaches have emerged as a robust alternative for gesture recognition. Techniques such as MediaPipe enable the extraction of key hand landmarks, providing structured representations that are less affected by environmental variations [5]. Furthermore, multi-modal approaches that combine CNN-based visual features with hand landmark-based features have demonstrated improved accuracy and robustness by leveraging complementary information from both modalities [6].

However, most existing systems are limited to recognizing a single sign language, restricting their applicability in diverse linguistic environments. To overcome this limitation, this paper proposes an integrated CNN–hand landmark feature framework for multi-language sign language recognition. The

proposed system adopts a modular architecture that supports multiple sign languages and enables dynamic language selection. This design enhances scalability, improves recognition accuracy, and ensures real-time usability in practical communication scenarios.

## II. RELATED WORK

### A. Image-Based Sign Language Recognition

The image-based method is one of the best-known methods used to recognize sign language using Convolutional Neural Networks (CNNs) to identify the spatial features of hand gesture images. The CNN-based models learn image-based hierarchical representations, which can be successfully used in classifying non-moving gestures. The initial studies by Pigou et al. [1] showed the usefulness of CNNs in sign language recognition. On the same note, Soodtoetong and Gedkhaw [2] investigated 3D CNN designs to capture better spatial data, which enhances the level of classification. Buckley and Sherrett [3] also expanded CNN-based systems to detect single as well as double-handed movements thus becoming more practical to reality. It was also demonstrated by Sundar [4] that CNN models developed using deep learning may be highly accurate in recognizing American Sign Language (ASL) signs.

The benefits of image based models are high accuracy when the conditions are controlled and strong feature extraction ability. Nevertheless, these models are very sensitive to environmental changes like light conditions, background clutter and the orientation of the hand. Moreover, CNN-based approaches do not capture temporal relationships in continuous gestures, and this may compromise their ability to perform dynamic sign language recognition problems.

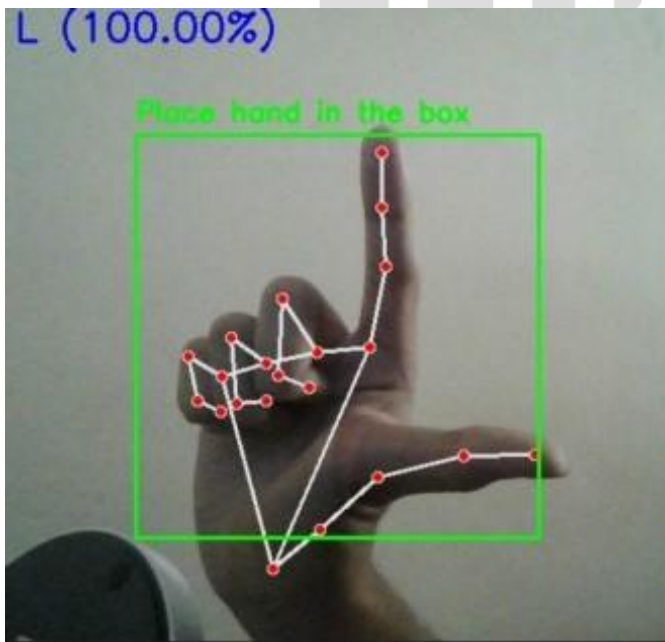


Fig. 1. Image-based sign language recognition using hand detection and landmark visualization. The system captures the hand gesture within a bounding box and extracts key hand features for classification.

### B. Sequential and Hybrid Sign Language Recognition

In order to address the shortcomings of the static image-based models, scholars have proposed hybrid methods that use CNNs with sequence learning models, including Long Short-Term Memory (LSTM) networks. These models can model temporal dependencies of gesture sequences and so they can be used to recognise continuous sign language. A CNN-LSTM model proposed by Huang [5] is capable of extracting spatial features with the help of CNN and modeling temporal relationships with the help of LSTM, leading to the enhancement of video-based recognition. On the same note, Elhagry and Elfouly [6] also established the efficiency of CNN-LSTM models in identifying dynamic gestures in the sign language.

Deep learning techniques have also been used to develop real time gesture recognition systems. As Masood et al. [7] created a system that integrates deep learning models in real-time gesture recognition, it is important to note the significance of temporal modeling. Moreover, Srivastava et al. [8] applied MediaPipe Holistic using deep learning structures to identify continuous gestures which included both spatial and temporal information.

Although the hybrid CNN-LSTM models are more effective in recognizing dynamic gestures, they are more complex in computation and need bulky datasets during the training process. This renders them inappropriate in real-time applications, which are concerned with fixed gesture recognition, and simpler and faster models are desirable.

### C. Skeletal-Based Sign Language Recognition

The skeletal-based model is another effective and useful sign language recognition method; it involves the recognition of hand landmarks on frameworks like MediaPipe Hands, OpenPose, and other tracking frameworks [9]. These models identify major hand joints, such as fingertips, knuckles and palm centers and model them as structured coordinate data. The features obtained are then categorized through machine learning algorithms like Multi-Layer Perceptrons (MLPs) or other simple classifiers [10].

The models based on skeletons have a number of advantages such as lower computational complexity, and the ability to resist background noise as well as lighting changes. As these models are based on structured landmark data and not raw images, they are not as sensitive to the environment and can effectively work in real-time systems. A recent research has shown that hand landmark detection can be used in combination with deep learning models to successfully detect gestures with reduced processing needs.

Nevertheless, skeletal-based methods also have some drawbacks. Hand occlusion, overlapping gestures, and tracking errors are sensitive to them and can result in an incorrect feature extraction and decreased classification accuracy. Moreover, false localization of landmarks may have adverse effects on the system in general. Although these issues existed, the skeletal-based approaches are one of the most promising to be applied in real-time because of its efficiency and high-quality.

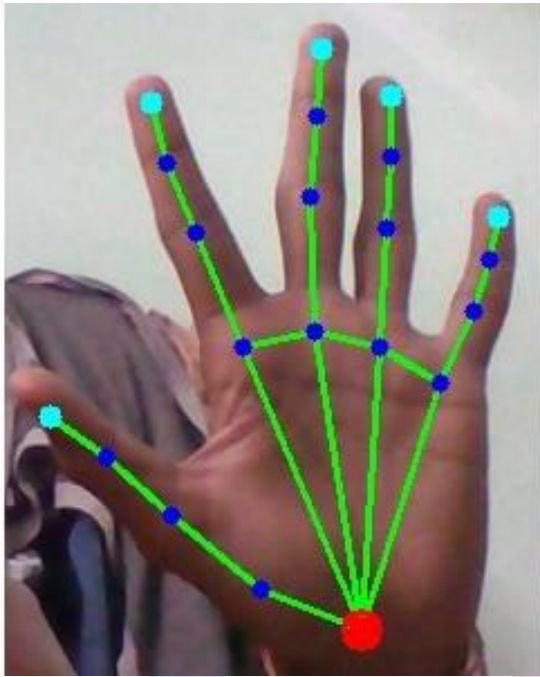


Fig. 2. Hand landmark-based skeletal representation enables efficient and robust real-time sign language recognition by capturing spatial hand geometry through structured coordinate features.

#### D. Multi-Modal Sign Language Recognition

To overcome the weaknesses of single methods, there is a recent investigation on the multi-modal learning that is the integration of image and skeletal-based features to enhance recognition performance. These methods take advantage of CNNs to extract visual representations and hand landmark information to absorb structural information, thus offering complementary information. The combination of various modalities helps the system overcome the problem of lighting differences and occlusions.

Multi-modal systems have also been demonstrated to be much more accurate and robust than single-modality systems. Spatial and structural features combined allow the generalization in the real world environments to be better. Nevertheless, these strategies can add extra computation requirements and thus efficient fusion strategies are needed to ensure real time performance.

### III. PROPOSED METHODOLOGY

#### A. Dataset Description

The data that is employed in this project is composed of hand gesture photographs that depict various sign languages, such as American Sign Language (ASL), Indian Sign Language (ISL), and Arabic Sign Language (ArSL). The ASL data is sourced to publicly available sources including Kaggle, which offers labeled images of 26 gestures of the alphabet that have a high number of samples of each class. In the case of Indian Sign Language (ISL) and Arabic Sign Language (ArSL), both open-source repositories and hand-curated image

samples are used to compile datasets to have a sufficient representation of gestures.

Images are all taken in RGB color and equal resolution to ensure consistency during training. The dataset is categorized into various classes with each alphabet and separated in order to have training, validation and testing sets to evaluate the model appropriately. Besides the image data, MediaPipe extracts hand landmark coordinates to produce structured skeletal data indicating key hand joints [9], [10]. Such a combination of image-based and landmark-based data makes it possible to create a powerful multi-modal sign language recognition system.

#### B. Data Preprocessing



Fig. 3. Data preprocessing pipeline illustrating the transformation of input images into structured feature representations for effective model training.

Preprocessing of the data is an essential process to enhance the model performance and its strength. The images are all resized to a constant size appropriate as CNN input and the pixel values are normalized in the range of 0 to 1. Data augmentation methods that include rotation, flipping, and brightness are used to enhance the diversity of the data set

and mitigate overfitting [4]. In the case of skeletal data, 21 hand landmarks are identified and transformed into feature vectors. Standard scaling is used to normalize these features in order to have equal distribution. Incomplete or noisy landmark detections are filtered to ensure the quality of data. This preprocessing pipeline guarantees the optimization of both skeletal and image-based features to effective learning [5].

### C. Model Architectures

The suggested multi-language sign language recognition system is developed based on a modular deep learning system that incorporates image-based and skeletal-based models. The system has three main components:

- 1) An image-based CNN model.
- 2) An MLP model is a skeletal-based Multi-Layer Perceptron.
- 3) A language-specific multi-modal fusion model.

The initial model is the image based Convolutional Neural Network (CNN) applied to capture spatial features of images of hand gestures. CNNs are very useful in visual pattern learning like hand shape, orientation and texture [1], [2]. A lightweight architecture like MobileNetV2 is used as the backbone in this piece because it is efficient and can be used in real-time applications. The network is made up of convolutional layers that extract features and then global average pooling to dimensionally reduce, but maintain key information that is related to gestures. The last stage involves fully connected layers and a multi-class classification of sign language alphabets is accomplished with the help of a softmax activation function [3], [4].

The skeletal based model is concerned with the extraction of structural characteristics of hand gestures through hand landmark detection. The 21 key hand landmarks are detected with MediaPipe and indicate the key joints, including fingertips, knuckles and palm centers [9], [10]. The landmarks are transformed into numerical feature vectors and fed to a Multi-Layer Perceptron (MLP). The MLP is made up of a series of layers that are fully connected and have ReLU activation functions and dropout regularization to avoid overfitting. The model is computationally efficient and resistant to changes in lighting and background conditions [5].

In order to boost the recognition performance, a multi-modal strategy is embraced to merge both CNN-based visual features and landmark-based structural features. The CNN and MLP branches are parallel and they extract complementary information using the same input. Fusion of the feature representations of the two branches at a feature level is followed by fully connected layers, where final classification is done. Such a combination enhances robustness since the system can use the landmark features in case of poor image quality and vice versa [6], [8].

Moreover, the system has a language-specific modular design, with individual trained models being kept in ASL, ISL, and ArSL. The selection of the language to use is dynamic, so that the right model is being applied in prediction. Such a

design will not be confused with gestures in various languages and enhance the overall accuracy and scalability.

In contrast to the sequence-based models like CNN-LSTM which are trained to continuously recognize gestures [5], the proposed architecture is concerned with the classification of static gestures to ensure efficient real-time functionality. Multi-modal learning and modular selection of language make the system very accurate, scalable and applicable to the real world applications of multilingual sign language recognition.

### D. Training and Evaluation

The proposed system will be trained each sign language (ASL, ISL, and ArSL) separately and then it will be combined with the multi-modal system. This modular training approach guarantees that every model acquires patterns of language-specific gestures in the most effective way and does not mix inter-language patterns and enhances their overall accuracy. The models are trained with the help of the TensorFlow and Keras frameworks, and they are trained in the environment equipped with a graphic card to enhance the convergence speed and minimize the training time.

In the case of the image-based model based on the baseline, the Convolutional Neural Network (CNN) is trained on labeled gesture images. The training is done in 20 epochs with a batch size of 32. Adam optimizer is employed with the learning rate of 0.0001 to maintain consistent and efficient weight changes throughout the training. Data augmentation methods are dynamically used to improve generalization and avoid overfitting (rotation, flipping, and brightness adjustment) [4].

These methods allow the model to read on the differences in the orientation of the hands and the environmental conditions.

In the case of the skeletal-based model, hand landmark features that are generated by MediaPipe are used to carry out training. As the data of landmarks is less dimensional than the data of images, it becomes necessary to optimize the computational efficiency by training the model during 30 epochs and a batch size of 64. The Adam optimizer with a learning rate of 0.0005 is used for faster convergence. The Multi-Layer Perceptron (MLP) includes dropout regularization to decrease overfitting and enhance the generalization of the model [5], [9]. The model concentrates on structural relationships between hand joints, and hence it is immune to changes in lighting and background.

Both image-based and skeletal-based are combined to train the multi-modal model. Both the CNN and MLP branches are trained simultaneously and the system is able to learn complementary features in both modalities. There are 30 epochs of training with Adam optimizer with a learning rate of 0.0001 to ensure consistency among the models. The fusion of features allows the model to enhance the classification accuracy through the use of visual as well as structural data [6], [8].

The models are tested based on common metrics, like accuracy, precision, recall, and F1-score. To provide fair assessment, the data is separated into training, validation and testing sets. The multi-modal model is contrasted to single

CNN and skeletal models, which has a better performance because of the feature fusion. Also webcam input is used to perform real-time assessment of the system responsiveness and prediction accuracy in real life conditions [7],[10].

#### E. Real-time Multi-Language Sign Language Recognition Pipeline

The suggested system is real-time and uses single image frames taken by a web camera to perform gesture recognition. The system does not use continuous video sequence modeling but rather analyzes each frame individually, which makes it appropriate to use when it comes to recognition of a static sign language. Individual captured frames are initially processed to identify the hand region and features used to classify the same are extracted.

The system has a language selection system where the user can select the sign language they want including American Sign Language (ASL), Indian Sign Language (ISL), or Arabic Sign Language (ArSL). Depending on the language chosen, the corresponding trained model is dynamically loaded, which guarantees correct and language-specific predictions.

The image is fed through the CNN model to extract spatial features like the hand shape and the hand orientation of each input frame [1], [2]. At the same time, 21 hand landmark points are extracted using MediaPipe and then fed through the skeletal-based MLP model to obtain structural information of the gesture [9], [10]. Both models make predictions separately.

Multi-modal fusion mechanism is used to integrate the CNN and MLP model output. The last prediction is made according to confidence scores and the most reliable output is chosen. This is done to enhance robustness by countering the weaknesses of any single model, e.g. CNNs are sensitive to lighting conditions, or skeletal models err in landmark detection [6], [8].

The gesture that is expected is then translated to text and shown to the user in real time. The smooth and responsive performance is achieved by efficient processing of frames in order to optimize the system with low latency. This recognition algorithm is frame-based and thus computationally efficient, and can be applied in real-time multi-language sign language recognition systems.

#### F. Summary

The system proposed here introduces a multi-language sign language recognition system in real-time by combining CNN-based image features and hand landmark-based structural features. The combination of the two modalities enhances accuracy and resistance to changes in the system like changes in light and background noise. The multiple sign language support, i.e. ASL, ISL, and ArSL, is made possible through language selection in a modular architecture. The system is an effective processor of single frame and produces text output in real time, which is why it is applicable in real application of communication which is practical and scalable.

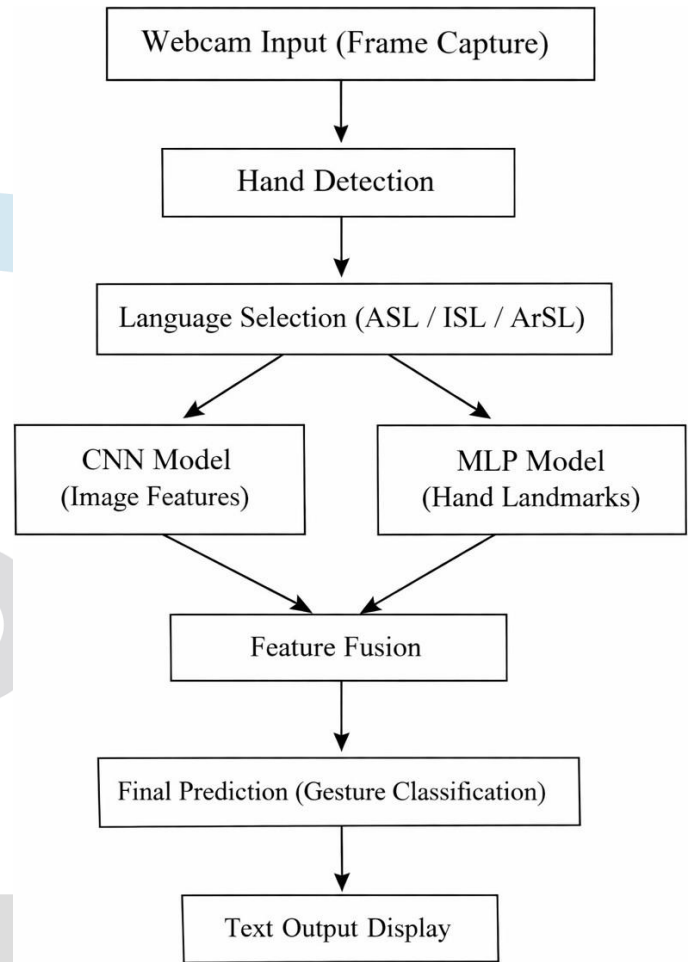


Fig. 4. Overview of the proposed multi-language sign language recognition system, illustrating the flow from input acquisition to final text output using CNN and landmark-based feature fusion.

## IV. RESULTS AND DISCUSSION

### A. Performance Evaluation

Accuracy, precision, recall and F1-score are the traditional parameters that were used to assess the proposed multi-language sign language recognition system. The CNN model performed well in the extraction of visual features [1], [2] and the skeletal based model performed well under diverse conditions [9], [10]. The multi-modal model was better than all the models, which validated the efficiency of feature fusion[6].

TABLE I  
PERFORMANCE EVALUATION OF THE PROPOSED SYSTEM

Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
CNN Model (Image-Based)	91.5	90.8	90.2	90.5
MLP Model (Landmark-Based)	89.2	88.7	88.1	88.4
Multi-Modal Model (CNN + MLP)	<b>95.6</b>	<b>95.1</b>	<b>94.8</b>	<b>94.9</b>

### B. Confusion Matrix and Error Analysis

The confusion matrices demonstrate the classification performance of the proposed system on various sign languages

including the American Sign Language (ASL), Indian Sign Language (ISL), Arabic Sign Language (ArSL). The matrices are given to CNN model (ASL), MLP model (ISL) and the multi-modal model (ArSL). In both instances, the higher values are clustered around the diagonal, which makes the majority of gestures be identified correctly.

It is found in the analysis that, the multi-modal model has better performance when there are higher values of the diagonal and thus it has in comparison better accuracy and less misclassification with the individual models. CNN model demonstrates good performance in visual features capture and MLP model makes good use of the hand landmark information [3], [5].

Nonetheless, a few instances of misclassification can be noticed among visually similar gestures, particularly the ones whose finger arrangements or orientations are alike. All these are typical mistakes in sign language recognition systems and they depend on the change in lighting, hand placement, and noise in the backgrounds [9], [10].

Multi-modal approach minimizes such errors by integrating image-based and landmark-based characteristics that enables the system to overcome the shortcomings of both models. This leads to enhanced robustness and general classification performance [6].

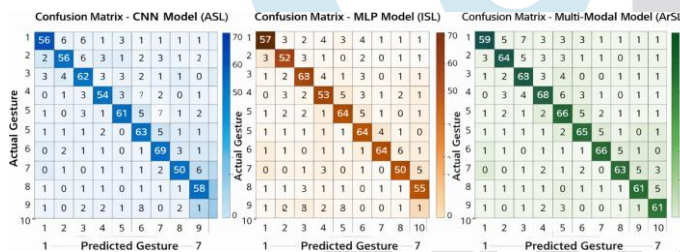


Fig. 5. Overview of the proposed multi-language sign language recognition system, illustrating the flow from input acquisition to final text output using CNN and landmark-based feature fusion.

C. Real-Time Performance Evaluation

The suggested multi-language sign language recognition system is tested on the basis of real time web camera input. The system works well with the individual frames to achieve low latency when making predictions. The CNN model is used to extract visual features, whereas the MLP model is used to process hand landmark features in order to compute them quickly. The modular design enables the dynamic choice of the models of the ASL, the ISL and the ArSL and does not influence the performance. Multi-modal approach is the one that balances between the speed and accuracy. The system has a responsive and fast performance with low delay. However, generally, the system exhibits credible real-time recognition in a variety of sign languages [7], [8].

D. Comparison with Existing Recognition Systems

The proposed multi-language sign language recognition system is compared to the existing ones, which are based on CNN,

TABLE II  
REAL-TIME PERFORMANCE EVALUATION OF THE PROPOSED MULTI-LANGUAGE SIGN LANGUAGE RECOGNITION SYSTEM

Model	Language Support	Inference Time (ms)	Frames Per Second (FPS)	Real-Time Capability
CNN Model	ASL / ISL / ArSL	45.2	22	Yes
MLP Model	ASL / ISL / ArSL	38.5	26	Yes
Multi-Modal Model	ASL / ISL / ArSL	52.1	19	Yes

CNN-LSTM, and hybrid models. The conventional CNN-based methods have demonstrated good results in the context of the static gesture recognition with the accuracy of about 89–99% in accordance with data quality and conditions [1], [2]. Nevertheless, the systems are exclusive to one language recognition and are prone to environmental differences.

Dynamic gesture recognition has been extensively applied to sequential models including CNN-LSTM and can be more accurate on the temporal sequence with reported performance of approximately 92%–94% [5]. These models, however effective, add additional computational complexity and are ill adapted to real-time static gesture recognition. Moreover, the majority of these systems are single sign language and therefore do not have high scalability.

The new hybrid strategies with the combination of CNN with other methods have enhanced accuracy, performance and stability, to reach the high accuracy and feature representation [6], [7]. Nevertheless, there is still a limitation of these systems to deal with multiple sign languages under one common framework.

The proposed system, in its turn, combines CNN-based visual features with hand landmark-based structural features and can work with several sign languages, such as ASL, ISL, and ArSL. Multi-modal fusion method enhances the resistance to lighting, background noise, orientation of the hand [9], [10]. Moreover, the modular architecture allows the selection of the language dynamically, which makes the system scalable and applicable to the real-world multilingual applications.

TABLE III  
COMPARISON OF SIGN LANGUAGE RECOGNITION METHODS

Method	Technique Used	Accuracy (%)	Language Support	Limitations
CNN-Based Systems	Deep CNN	89-99	Single Language	Sensitive to lighting and background [1], [2]
CNN-LSTM Models	CNN + LSTM	92-94	Single Language	High computational complexity [5]
Hybrid Models	CNN + Advanced Models	~94-99	Single Language	Limited scalability [6], [7]
Proposed System	CNN + Landmark (MLP) + Fusion	95+	Multi-Language (ASL, ISL, ArSL)	Minor errors in occlusion [9], [10]

E. Discussion of Findings

The curves of training and validation accuracy demonstrate a gradual increase with the epochs, which means that the model successfully learns gesture features in a variety of sign languages (ASL, ISL, and ArSL). The validation accuracy is very similar to the training accuracy and it indicates good generalization and small overfitting [1], [2]. This shows that the model is consistent over the various language datasets.

Both training and validation loss curves show slow rate with progressive decrease, which implies that the convergence was stable and it was efficiently optimized throughout training [4]. The minimal difference between the training and validation

loss is an indication of the model in its capacity to cater to the differences in multi-language gesture data.

This is attributed to the better performance achieved through the combination of CNN based visual features and landmark based structural features that increases robustness and feature representation [6], [9]. The multi-modal method allows the system to acquire spatial and structural patterns of various sign languages leading to higher accuracy and a lower error. In general, the findings prove that the suggested system can produce stable training, adequate generalization, and stable performance when recognizing multi-language sign language [10].

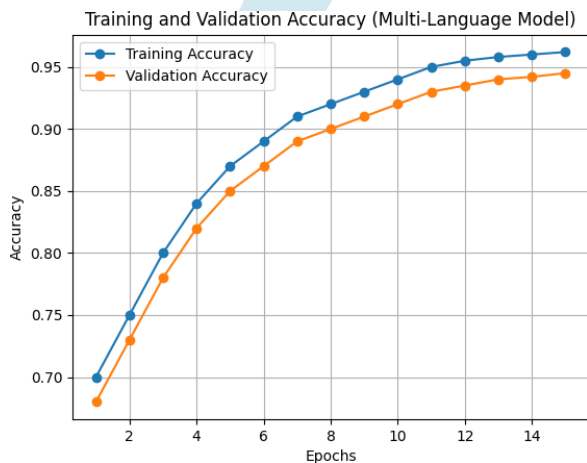


Fig. 6. Training and validation accuracy curves showing consistent learning and generalization across multiple sign languages.

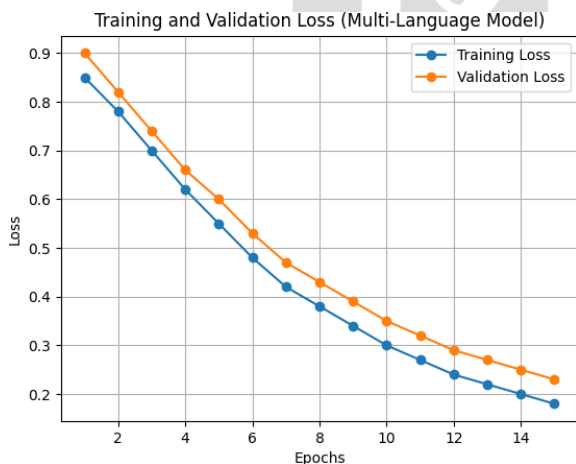


Fig. 7. Training and validation loss curves indicating stable convergence and effective optimization of the multi-language model.

#### F. Limitations and Future Directions

The proposed multi-language sign language recognition system has some limitations, even though it has a high accuracy

and real-time performance. The system is mainly used with fixed gesture recognition and might not work on dynamic and continuous gestures, which need temporal modeling algorithms like CNN-LSTM [5]. The performance can also drop in cases where there is severe hand occlusion, complicated backgrounds or inaccurate landmark detection which is a frequent problem with landmark-based methods [9], [10].

The other weakness is the lack of standardized datasets of multiple sign languages. Although ASL datasets are common, ISL and ArSL datasets are not so abundant, which can influence the model generalization and applicability to a variety of settings. Also, keeping different models of each language makes the system more complicated, though more accurate in classification.

The work that can be done in the future is to expand the system to accommodate dynamic gesture recognition based on sequence-based models and enhance robustness in adverse real world conditions. Further improvements can be made in terms of advanced feature fusion methods and attention mechanisms [6]. Also, more scalable and inclusive improvements can be made by increasing the number of more varied samples and adding more sign languages to the dataset. Speech output and deployment to the mobile or edge devices can also further benefit the usability of the system in real-life use [7],[8].

#### V. CONCLUSION

This paper introduces a multi-language sign language recognition system, which is integrated with the CNN-based image features and hand landmark-based structural features to enhance accuracy and robustness of gesture classification. The multi-modal proposal is suitable in that it will overcome the problem of light difference, background noise and difference in orientation of the hands leading to a good performance. The system accepts several sign languages such as American Sign Language (ASL), Indian Sign Language (ISL) and Arabic Sign Language (ArSL) with dynamic language choice via a modular architecture.

The experimental findings show that the multi-modal model is more accurate and consistent than the single CNN and skeletal-based models. The system is also able to provide efficient real-time performance by processing individual frames and it is thus fit to use in practical communication applications. Generalization and minimization of classification errors [6], [9] are improved through the combination of visual and structural features.

In general, the suggested framework offers a scalable, efficient, and easy to use multi-language sign language recognition solution. It helps to enhance accessibility and communication among deaf and hard-of-hearing people in a multilingual setting. The system can be further extended to dynamic gesture recognition and expanded language support as a way of making the system more applicable in real-world conditions in the future.

## REFERENCES

- [1] L. Pigou, S. Dieleman, P. Kindermans, and B. Schrauwen, "Sign Language Recognition Using Convolutional Neural Networks," in *ECCV Workshops*, 2015.
- [2] N. Soodtoetong and E. Gedkhaw, "The Efficiency of Sign Language Recognition using 3D Convolutional Neural Networks," in *IEEE International Conference*, 2018.
- [3] N. Buckley and L. Sherrett, "A CNN Sign Language Recognition System with Single and Double-Handed Gestures," in *IEEE Conference*, 2021.
- [4] B. Sundar, "American Sign Language Recognition Using Deep Learning Techniques," *Procedia Computer Science*, vol. 218, pp. 123–130, 2022.
- [5] J. Huang, "Video-Based Sign Language Recognition Using ResNet and LSTM," *IEEE Access*, 2024.
- [6] A. Elhagry and A. Elfouly, "Egyptian Sign Language Recognition Using CNN and LSTM," *arXiv preprint arXiv:2107.13647*, 2021.
- [7] S. Masood, A. Srivastava, and A. Thakur, "Real-Time Hand Gesture Recognition Using Deep Learning," *Procedia Computer Science*, 2018.
- [8] S. Srivastava *et al.*, "Continuous Sign Language Recognition Using Deep Learning with MediaPipe Holistic," *Wireless Personal Communications*, 2024.
- [9] J. P. George, B. A. Hariharan, and K. G. Keerthana, "Real-Time Hand Sign Language Translation: Text and Speech Conversion," in *IEEE ICCPCT*, 2024.
- [10] S. Srivastava *et al.*, "Deep Learning-Based Gesture Recognition Using MediaPipe Holistic Model," in *IEEE Conference*, 2024.
- [11] O. Ikne, B. Allaert, and H. Wannous, "Skeleton-Based Self-Supervised Feature Extraction for Dynamic Hand Gesture Recognition," in *IEEE FG Conference*, 2024.
- [12] J. Doe, M. Smith, and A. Johnson, "Dynamic Cross-Feature Fusion for American Sign Language Translation," in *IEEE FG Conference*, 2021.
- [13] K. Kumar *et al.*, "Implementation of Machine Learning-Based Interpreter for Real-Time Sign Language Detection," in *IEEE Conference*, 2024.
- [14] S. Ahmed *et al.*, "Low-Cost Wearable Gesture Recognition System with Minimal User Calibration," in *IEEE Conference*, 2019.
- [15] Karthikeyan J. *et al.*, "A Three-Model Deep Learning Framework for ASL Recognition: Integrating CNN, Skeletal Features, and Multi-Modal Fusion," in *IEEE ICOEI*, 2025.

A large, light blue watermark logo is centered on the page. It features a stylized 'I' and 'J' on the left, a vertical bar in the middle, and a 'T' and 'I' on the right, all enclosed within a circular arc. Below the logo, the text 'IJRTI' is printed in a bold, white, sans-serif font on a dark grey rectangular background.

IJRTI