# DESIGNING A MULTIMODAL SIGN LANGUAGE RECOGNITION SYSTEM WITH CNN-RNN FUSION

## ABSTRACT

Sign language remains a crucial means of communication for the deaf and hard-of-hearing population, yet accessible, real-time translation solutions for Indian Sign Language (ISL) are still limited. An advanced recognition system has been developed to address this gap, leveraging both hand skeleton and facial expression features to enhance recognition accuracy and robustness. The MediaPipe framework enables precise extraction of hand and facial landmarks from video frames, capturing the subtle gestures and expressions essential for accurate interpretation. Landmark sequences are processed using a hybrid deep learning architecture that integrates a one-dimensional Convolutional Neural Network (1D CNN) for spatial feature extraction with a Recurrent Neural Network (RNN) for temporal sequence modeling. Training and evaluation on the ISL and ISL-CSLTR datasets, which encompass a wide variety of isolated and continuous sign language samples, demonstrate the system's effectiveness. Experimental results, as illustrated by the training and validation accuracy/loss curves and a highly diagonal confusion matrix, indicate high accuracy in both word-level and sentence-level recognition tasks, exceeding the performance of traditional methods that rely solely on manual features. Real-time operation further supports practical deployment in communication aids and assistive technologies. By combining hand and facial cues with state-of-the-art deep learning, the solution offers a robust and scalable approach to automated ISL recognition, with significant potential to bridge communication barriers and foster greater inclusivity for the deaf community in India.

# CHAPTER 1

# INTRODUCTION

## 1.1  ABOUT THE PROJECT

Sign language serves as an essential communication medium for millions of deaf and hard-of-hearing individuals worldwide. Unlike spoken languages, sign language relies on a combination of hand gestures, body movements, and facial expressions to convey meaning. In India, Indian Sign Language (ISL) is widely used, yet a significant communication barrier persists between sign language users and the broader hearing population. Automated sign language recognition (SLR) systems have the potential to bridge this gap, enabling seamless interaction and fostering inclusivity. Traditional SLR approaches often struggle with the complexity and variability of gestures, especially when both manual (hand) and non-manual (facial) cues are involved. Recent advancements in computer vision and deep learning, particularly the use of frameworks like MediaPipe for landmark detection and hybrid neural network architectures, have opened new avenues for more accurate and robust SLR systems.

## 1.2  OBJECTIVE OF THE PROJECT

The project focuses on designing and implementing a real-time Indian Sign Language recognition system that integrates both hand skeleton and facial expression features to enhance accuracy and reliability. A robust feature extraction pipeline is developed using MediaPipe, enabling precise detection of hand and facial landmarks from video input. These extracted features serve as the foundation for constructing a hybrid deep learning model that combines a one-dimensional Convolutional Neural Network (1D CNN) with a Recurrent Neural Network (RNN). This architecture is adept at learning spatial patterns from individual frames and temporal dependencies across gesture sequences, which are essential for accurate sign language interpretation. The system

undergoes comprehensive training and evaluation using the ISL and ISL- CSLTR datasets, which collectively provide a wide spectrum of isolated signs and continuous sign language sentences. Through this approach, the model attains high recognition accuracy and demonstrates real-time performance, making it suitable for practical deployment in diverse real-world applications where accessible and efficient sign language recognition is required.

## 1.3  SCOPE OF THE PROJECT

The project is centered on recognizing Indian Sign Language at both word and sentence levels by leveraging video data. A real-time sign language recognition system has been implemented, capable of processing both live video streams and pre-recorded videos for versatile application. To enhance recognition accuracy, the system extracts and utilizes both hand and facial landmarks, effectively addressing the limitations found in approaches that rely solely on hand gesture analysis. Advanced deep learning techniques,  specifically a hybrid 1D CNN combined with RNN, are employed to achieve effective spatial-temporal feature learning. System performance is rigorously evaluated using benchmark datasets such as ISL and ISL-CSLTR, ensuring comprehensive coverage across a wide range of vocabulary and sentence structures. While speech synthesis and translation to other spoken languages are not included in the current implementation, the architecture is designed to support future integration of such features. Scalability has been prioritized, enabling potential expansion to additional sign languages and facilitating deployment on a variety of platforms, including mobile devices.

# CHAPTER 2
# LITERATURE SURVEY

**Paper 1:**
**Title: "Research Paper on Real-Time Sign Language Interpreter using Mediapipe Holistic"**
**Authors:** S.R. Aajmane, A.S. Neje, B.S. Khedkar, S.S. Koulage, S.M. Momin
**Year:** 2022

The paper presents a real-time sign language interpreter leveraging the MediaPipe Holistic framework to bridge communication gaps for the deaf and hard-of-hearing community. The system processes video input, detecting and plotting hand gestures, body posture, and facial expressions using MediaPipe Holistic's landmark extraction. These extracted keypoints are then fed into an AI model for sign recognition, with additional natural language processing (NLP) techniques applied to improve sentence-level interpretation. The workflow involves segmenting video frames, tracking hand movements via the CamShift method, and employing the P2DHMM algorithm for robust hand tracking. MediaPipe Holistic enables the system to combine pose, hand, and face landmarks into a unified representation, which is crucial for accurately capturing the nuances of sign language. The paper highlights the system's ability to interpret signs in real-time, converting them into meaningful sentences rather than isolated words, which enhances communication effectiveness for users. The authors note that integrating MediaPipe Holistic with deep learning models (such as LSTM) allows for efficient processing and high accuracy, even with relatively small datasets. The approach demonstrates the potential of combining computer vision and NLP for practical, real-time sign language translation, though challenges remain in scaling to large vocabularies and diverse sign languages.

**Paper 2:**

**Title:** "Egyptian Sign Language Recognition Using CNN and LSTM"

**Authors:** Ahmed Adel Gomaa Elhagry, Rawan Gla Elrayes, Dr.Amr Zamel and Dr.Ahmed Helmy

**Year:** 2021

The paper presents a video-based Egyptian Sign Language (ESL) recognition system designed to support the deaf community in Egypt. The authors explore two neural network architectures: a CNN-only model, which uses a retrained Inception v3 network to extract spatial features from video frames and achieves 90% accuracy for classifying isolated ESL signs, and a CNN-LSTM model, which combines CNN for spatial feature extraction with a Long Short-Term Memory (LSTM) network to capture temporal dynamics in signing, achieving 72% accuracy. The study highlights challenges such as regional variations in ESL and the limited availability of large, diverse datasets. The authors suggest expanding the dataset and exploring additional deep learning models for future work, aiming to develop a mobile application that facilitates real-time communication for the Egyptian deaf community.

**Paper 3:**

**Title:** "Indian Sign Language Recognition Using Mediapipe Holistic"

**Authors:** Kaushal Goyal and Dr. Velmathi G

**Year:** 2023

The paper presents a robust system for Indian Sign Language (ISL) recognition using the MediaPipe Holistic framework, aiming to convert ISL gestures into text or speech to bridge communication gaps for deaf and hard-of-hearing individuals. MediaPipe Holistic is leveraged to extract key landmarks from the hands, face, and body, significantly reducing dataset size and computational requirements compared to video-based approaches. The system compares Convolutional Neural Network (CNN) and Long Short-Term Memory

(LSTM) models: CNNs excel at recognizing static signs (letters and characters), while LSTMs outperform CNNs for dynamic gestures and sentences by tracking temporal movement of hands, face, and pose. The methodology focuses on extracting only essential keypoints, storing them efficiently, and enabling real-time detection. This approach avoids data-heavy image operations, making the model accessible and easy to train for new users. The authors highlight the scarcity of ISL-specific datasets and research, noting that most prior work focuses on other sign languages. Their system demonstrates improved recognition accuracy and practical usability for both static and dynamic ISL, with potential for further enhancement through expanded datasets and text-to- sign language translation models.

**Paper 4:**

**Title: "Real Time Indian Sign Language Recognition using Convolutional Neural Network"**

**Authors:** Snehal Hon, Manpreet Sidhu, Sandesh Marathe and Tushar A. Rane

**Year:** 2024

The paper presents a real-time Indian Sign Language (ISL) recognition system leveraging Convolutional Neural Networks (CNNs) to interpret hand gestures captured through a camera, aiming to bridge communication gaps for the deaf and hard-of-hearing community. The proposed approach uses a hybrid CNN framework, often enhanced with techniques like transfer learning (e.g., VGG16, ResNet152) and grid-based feature extraction, to achieve robust recognition of ISL gestures without the need for external hardware. Key steps include image preprocessing, gesture detection, and classification, with the CNN model trained on large datasets of static hand gesture images. The system achieves high accuracy-typically exceeding 90% and reaching up to 99.44% in some implementations-demonstrating its effectiveness for both static and dynamic ISL recognition in real-time scenarios. The models are designed to be

lightweight and deployable on mobile devices, using frameworks like TensorFlow Lite, making them accessible and practical for everyday use. The study underscores the importance of tailored preprocessing, diverse feature extraction, and thoughtful model architecture in achieving reliable sign language interpretation. The integration of CNNs with real-time processing enables efficient, accurate, and user-friendly ISL recognition, supporting social inclusion and accessibility for the deaf community.

**Paper 5:**

**Title: "Continuous Sign Language Recognition with Correlation Network"**
**Authors:** Lianyu Hu, Liqing Gao, Zekang Liu and Wei Feng
**Year:** 2023

The paper introduces the Natural Language-Assisted Sign Language Recognition (NLA-SLR) framework, addressing the challenge of visually indistinguishable signs (VISigns) in sign languages-signs that appear similar but have different meanings. Traditional vision-based neural networks struggle to differentiate these VISigns, limiting recognition accuracy. The NLA-SLR framework leverages semantic information from glosses (sign labels) to enhance recognition. It employs language-aware label smoothing, where for VISigns with similar meanings, the framework generates soft labels using semantic similarities between glosses, easing model training and improving generalization. For VISigns with distinct meanings, it utilizes inter-modality mixup, blending features from both visual (video/keypoints) and gloss modalities to maximize the separability of different signs. The backbone of the framework is the Video-Keypoint Network (VKNet), which processes both RGB video and human body keypoints, capturing rich visual and motion cues across different temporal windows. Experiments on benchmarks like MSASL, WLASL, and NMFs-CSL show that the NLA-SLR achieves state-of-the-art performance, particularly excelling at distinguishing VISigns where traditional

models fail. The approach introduces minimal computational overhead and demonstrates that integrating natural language semantics directly into sign language recognition pipelines can substantially boost accuracy and robustness.

**Paper 6:**

**Title: "Sign Language Recognition Based on Facial Expression and Hand Skeleton"**

**Authors:** Zhiyu Long, Xingyou Liu, Jiaqi Qiao and Zhi Li

**Year:** 2024

The paper addresses the limitations of traditional sign language recognition (SLR) systems that rely solely on monocular cameras, which often suffer from low accuracy and poor robustness due to ineffective feature extraction and sensitivity to hand position and angle variations. To overcome these challenges, the authors propose a novel SLR network that integrates both hand skeleton features and facial expression information. Key innovations of the proposed approach include a hand skeleton feature extraction method that uses a coordinate transformation technique to more accurately describe hand shapes and spatial posture, thereby reducing interference from different hand positions and angles. Additionally, the network incorporates facial expression features, extracted using a Self-Cure Network (SCN), to provide additional contextual information that is often ignored in other SLR methods. This integration of facial cues enhances recognition accuracy and robustness, especially for signs where facial expressions are essential. Furthermore, both hand skeleton and facial expression features are input together into a data compensation network, allowing the model to compensate for missing or ambiguous information in either modality. Experiments on the LSA64 (Argentinian Sign Language) and SEUCSLRD (Chinese Sign Language) datasets demonstrate that the multimodal approach significantly improves recognition performance, particularly under varying conditions and interference.

**Paper 7:**

**Title: "Optimizing Hand Region Detection in MediaPipe Holistic Full-Body Pose Estimation to Improve Accuracy and Avoid Downstream Errors"**

**Authors:** Amit Moryossef

**Year:** 2024

The paper investigates a key limitation in MediaPipe Holistic's hand Region of Interest (ROI) prediction, which is critical for accurate full-body pose estimation and downstream applications like sign language recognition. The authors identify that the current heuristic for hand ROI works well only when the hand is parallel to the camera, but fails with non-ideal hand orientations, leading to inaccurate hand keypoint detection and subsequent errors in pose estimation. To address this, the paper proposes a data-driven approach that enriches the feature set used for ROI estimation. Instead of relying solely on three hand keypoints (wrist, index, and pinky), the new method incorporates additional keypoints (including the thumb, shoulder, and elbow) and leverages the z-dimension for more robust 3D localization. The authors experimented with both Kolmogorov-Arnold Networks (KANs) for interpretability and standard multilayer perceptrons (MLPs) for practical deployment, ultimately using MLPs due to framework limitations. Their approach resulted in significantly improved ROI estimation, as measured by higher Intersection-over-Union (IoU) scores compared to the baseline MediaPipe method. The improvements directly translate to better hand keypoint detection, which is essential for downstream tasks such as sign language recognition. The code and trained models are made publicly available, facilitating reproducibility and further research. The paper's critical contribution lies in demonstrating that a richer and more data-driven ROI estimation pipeline can substantially reduce errors in hand detection, especially in challenging scenarios. However, the solution is tailored to MediaPipe Holistic, and future work could explore generalizing the approach to other pose estimation frameworks.

**Paper 8:**

**Title: "An Open-Source Gloss-Based Baseline for Spoken to Signed Language Translation"**

**Authors:** Amit Moryossef, Mathias Muller, Anne Gohring, Zifan Jiang, Yoav Goldberg and Sarah Ebling

**Year:** 2023

The paper introduces an open-source, reproducible baseline system for spoken-to-signed language translation, designed to make sign language translation research more accessible and comparable across studies. The approach uses a modular pipeline with three main components for text-to-gloss translation: a lemmatizer (reducing words to their base forms), a rule-based word reordering and dropping system, and a neural machine translation (NMT) model. For gloss-to-pose conversion, the system leverages lexicon-derived data for Swiss German, Swiss French, and Swiss Italian Sign Languages, extracting skeletal poses from videos. The final step stitches these pose sequences into sign language sentences, producing photorealistic video outputs. The open-source implementation, available as a Python application and HTTP endpoint, enables community collaboration, reproducibility, and adaptation to resource-limited settings. This baseline is significant because it provides a common reference for evaluating and improving sign language translation systems, addressing the lack of standardization in the field.

**Paper 9:**

**Title: "Sign language recognition using the fusion of image and hand landmarks through multi-headed convolutional neural network"**

**Authors:** Refat Khan Pathan, Munmun Biswas, SuraiyaYasmin, Mayeen Uddin Khandaker, Mohammad Salman & AhmedA. F.Youssef

**Year:** 2023

The study proposes a cost-effective technique for American Sign

Language (ASL) recognition by fusing raw image data with hand landmark features using a multi-headed convolutional neural network (CNN). The method processes images in two parallel streams: one head analyzes the entire image, while the other extracts and processes hand landmarks-keypoints detected using a pre-trained model, similar to transfer learning. Data augmentation and dynamic learning rate reduction are used to prevent overfitting. Tested on the ASL Fingerspelling dataset (24 letters, complex backgrounds), the model achieved a high-test accuracy of 98.98%, demonstrating robustness in real- world scenarios where backgrounds and lighting vary. The dual-input approach- combining global image context and precise hand pose information-proved more effective than using either input alone. Limitations include reliance on the quality of hand landmark extraction and the need for a sufficiently large dataset for optimal training.

**Paper 10:**

**Title: "Siformer: Feature-isolated Transformer for Efficient Skeleton-based Sign Language Recognition"**
**Authors:** Muxin Pu, Mei Kuan Lim and Chun Yong Chong
**Year:** 2025

The paper presents Siformer, a novel transformer-based architecture designed for efficient skeleton-based sign language recognition (SLR). Siformer addresses three major limitations in existing skeleton-based SLR methods: reliance on unrealistic hand skeletal representations, vulnerability to missing data, and inefficient inference that does not account for varying sign complexity. The key innovations of Siformer include kinematic hand pose rectification, which enforces anatomical constraints on hand joints to produce more realistic and informative skeletal data for recognition. Additionally, the feature-isolated mechanism processes local spatial-temporal features independently, making the model robust to missing or incomplete joint data and

enhancing recognition accuracy in real-world scenarios. Furthermore, the input-adaptive inference component dynamically adjusts computational effort based on the complexity of the input sign, improving efficiency by allowing early stopping for simpler signs. Siformer achieves state-of-the-art results on benchmark datasets, with a top-1 accuracy of 86.50% on WLASL100 and 99.84% on LSA64, outperforming previous methods. The model's lightweight design and robustness to missing data make it suitable for practical deployment, including on mobile devices.


**Paper 11:**

**Title: "Emotion-Aware Indian Sign Language Recognition: A Multimodal Approach with Sign-Expression Correlation Analysis"**
**Authors:** Hriday Ranka, Darshit Sarda, Haardhik Kunder, Aaditya Rajesh and Aniket Kore
**Year:** 2024

The research introduces a real-time Indian Sign Language (ISL) recognition system that uniquely integrates both sign gesture interpretation and emotion detection to enhance communication for the deaf and hard-of-hearing community. The system employs Convolutional Neural Networks (CNNs) for facial emotion analysis and Long Short-Term Memory (LSTM) models for dynamic sign recognition, complemented by Natural Language Processing (NLP) for sentence construction. A major innovation is the creation of a custom ISL dataset, tailored to local linguistic and cultural nuances, which improves recognition accuracy over generic datasets. The emotion detection module analyzes facial expressions to provide additional context about the signer's emotional state, ensuring more natural and expressive communication. Sign-expression correlation analysis is used to verify that recognized signs and detected emotions are contextually consistent, further refining the system's reliability. The model achieves a high accuracy of 98.12% for sign recognition

and 79.22% for emotion detection on test data, demonstrating robust performance. Additionally, the system features text-to-speech conversion for accessibility, supporting real-world deployment in settings such as hospitals, schools, and public services.

**Paper 12:**

**Title: "Sign Language Recognition: A Comprehensive Review of Traditional and Deep Learning Approaches, Datasets, and Challenges"**
**Authors:** Tangfei Tao, Yizhe Zhao, Tianyu Liu and Jieli Zhu
**Year:** 2025

The comprehensive review surveys the evolution of sign language recognition (SLR), covering both traditional and deep learning approaches, the development of datasets, and ongoing challenges in the field. Traditional SLR methods initially relied on data gloves and sensor-based technologies, which offered high accuracy in capturing hand movements but were costly and lacked portability. Vision-based approaches using cameras became more popular due to their convenience, though they are more sensitive to environmental factors like lighting and background clutter. Early computer vision methods focused on handcrafted features and shallow classifiers, which struggled with variability in signer style and complex gestures. The advent of deep learning, particularly Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) such as Long Short-Term Memory (LSTM) networks, revolutionized SLR by enabling automatic feature extraction from raw images or video frames and modeling temporal dependencies in continuous signing. CNNs excel at static gesture recognition, while RNNs and LSTMs are adept at capturing the dynamic, sequential nature of sign language. Recent advances include attention mechanisms and multi-modal fusion, further boosting recognition accuracy. Dataset availability has been a key driver of progress. Large-scale datasets like MS-ASL and WLASL for American Sign Language, as well as multi-view

datasets for other languages, have enabled robust model training and benchmarking. Multi-view datasets, in particular, significantly improve recognition accuracy by providing diverse perspectives of the same sign, with studies showing up to 19.75% performance gains over single-view data. However, most datasets are collected in controlled environments and may not generalize well to real-world scenarios.

**Paper 13:**
**Title: "Natural Language-Assisted Sign Language Recognition"**
**Authors:** Ronglai Zuo, Fangyun Wei and Brian Mak
**Year:** 2023

The paper introduces the Natural Language-Assisted Sign Language Recognition (NLA-SLR) framework, addressing the challenge of visually indistinguishable signs (VISigns) in sign languages-signs that appear similar but have different meanings. Traditional vision-based neural networks struggle to differentiate these VISigns, limiting recognition accuracy. The NLA-SLR framework leverages semantic information from glosses (sign labels) to enhance recognition. It employs language-aware label smoothing, where for VISigns with similar meanings, the framework generates soft labels using semantic similarities between glosses, easing model training and improving generalization. For VISigns with distinct meanings, it utilizes inter-modality mixup, blending features from both visual (video/keypoints) and gloss modalities to maximize the separability of different signs. The backbone of the framework is the Video-Keypoint Network (VKNet), which processes both RGB video and human body keypoints, capturing rich visual and motion cues across different temporal windows. Experiments on benchmarks like MSASL, WLASL, and NMFs-CSL show that the NLA-SLR achieves state-of-the-art performance, particularly excelling at distinguishing VISigns where traditional models fail. The approach introduces minimal computational overhead and

demonstrates that integrating natural language semantics directly into sign language recognition pipelines can substantially boost accuracy and robustness.


**Paper 14:**

**Title: "Towards Online Continuous Sign Language Recognition and Translation"**

**Authors:** Ronglai Zuo, Fangyun Wei and Brian Mak

**Year:** 2024

The paper addresses the challenge of enabling real-time, online continuous sign language recognition (CSLR) and translation, which are essential for effective communication between deaf and hearing individuals. Traditional CSLR models typically rely on connectionist temporal classification (CTC) loss and require the entire video input for inference, causing high latency and memory usage, making them unsuitable for live or real-time scenarios. The proposed framework introduces a new approach with three main phases. First, a sign dictionary is developed, aligned with the glossary of the target dataset. Next, an isolated sign language recognition (ISLR) model is trained on this dictionary using both classification and saliency losses to improve accuracy. During inference, a sliding window method is used, where segments of the sign video are processed sequentially, enabling real-time prediction and reducing latency. For translation, the system incorporates a gloss-to-text network that uses the wait-k policy, allowing the translation process to begin as soon as glosses are recognized, rather than waiting for the entire sequence. Additionally, a lightweight post-processing step is introduced to remove duplicate predictions, and the framework can enhance existing offline CSLR models by combining features from both online and offline systems.

**Paper 15:**

**Title: "A Simple Baseline for Spoken Language to Sign Language Translation with 3D Avatars"**

**Authors:** Ronglai Zuo, Fangyun Wei, Zenggui Chen, Brian Mak, Jiaolong Yang and Xin Tong

**Year:** 2024

The paper introduces a practical system for translating spoken language into sign language, visualized through a 3D avatar-a task termed Spoken2Sign. Unlike earlier work that focused on sign-to-spoken translation, the approach enables hearing individuals to communicate with the deaf community in real time using sign language avatars. The system operates through a three-step pipeline. First, a gloss-video dictionary is created by segmenting continuous sign language videos from existing datasets into isolated signs using a state-of- the-art continuous sign language recognition model with CTC loss, mapping glosses (sign labels) to video clips. Next, each sign video is converted into a 3D pose representation, allowing the system to synthesize sign language from any viewpoint and enhancing the realism of the avatar's movements. Finally, the Spoken2Sign model is trained, comprising a Text2Gloss translator (which converts spoken language to glosses), a sign connector (which arranges glosses into sign sequences), and a rendering module that animates the 3D avatar using the gloss-3D sign dictionary. The output is a 3D avatar that signs the translated message, providing a flexible and visually rich communication channel. The paper also notes that 3D keypoint augmentation and multi-view understanding techniques can further enhance keypoint-based sign language understanding. This work is the first to present Spoken2Sign translation in a 3D avatar format, offering an open-source baseline and tools for further research and practical deployment in accessibility technologies.

# CHAPTER 3
# SYSTEM ANALYSIS

## 3.1 EXISTING SYSTEM

Sign language recognition (SLR) has been an active area of research for several decades, with early systems primarily relying on traditional image processing techniques and handcrafted features. In the context of Indian Sign Language (ISL), existing systems typically fall into two categories: sensor- based and vision-based approaches. Sensor-based systems utilize hardware such as data gloves or motion sensors to capture hand movements. While these systems can provide accurate data, they are often expensive, intrusive, and uncomfortable for users, limiting their practical adoption. Vision-based systems use cameras to capture gestures and apply image processing algorithms to interpret signs. Early vision-based SLR systems relied on background subtraction, skin color segmentation, and contour detection to isolate hand regions. Features such as Hu moments, Histogram of Oriented Gradients (HOG), and Scale-Invariant Feature Transform (SIFT) were commonly used to describe hand shapes and movements. These features were then classified using traditional machine learning algorithms like Support Vector Machines (SVM) or k-Nearest Neighbors (k-NN).

**Limitations:**

- Limited Feature Representation: Handcrafted features may not capture the full complexity and variability of sign language gestures, especially in the presence of occlusions, varying lighting conditions, and diverse backgrounds.

- Lack of Temporal Modeling: Many existing systems focus on static image recognition, neglecting the temporal dynamics of gestures that are crucial for sentence-level recognition.

- Neglect of Non-Manual Features: Facial expressions and head

movements, which are essential components of sign language, are often ignored, resulting in reduced recognition accuracy.

- Scalability Issues: Traditional methods struggle to scale to large vocabularies and continuous sign language recognition tasks.

Recent advancements in deep learning have led to the adoption of Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) for SLR. While these methods have improved performance, many still rely solely on RGB images or optical flow, without explicitly leveraging hand and facial landmarks. Furthermore, real-time performance and robustness in unconstrained environments remain challenging for many existing solutions.

## 3.2 PROPOSED SYSTEM

The proposed system aims to overcome the limitations of existing SLR solutions by integrating advanced computer vision and deep learning techniques for robust, real-time Indian Sign Language recognition. The system is designed to recognize both isolated signs and continuous sentences by leveraging both hand skeleton and facial expression features.

**Key Features of the Proposed System:**

- Landmark-Based Feature Extraction: Utilizes the MediaPipe framework to extract precise hand and facial landmarks from each video frame. This approach captures the spatial configuration of both manual (hand) and non-manual (facial) components, providing a rich feature set for recognition.

- Hybrid Deep Learning Model: Implements a combination of a one-dimensional Convolutional Neural Network (1D CNN) and a Recurrent Neural Network (RNN), specifically Long Short-Term Memory (LSTM) units. The 1D CNN extracts spatial features from sequential landmark vectors, while the RNN models temporal dependencies across frames, enabling accurate recognition of dynamic gestures and sentences.

- Comprehensive Dataset Utilization: Trains and evaluates the system using both the ISL dataset (for isolated word recognition) and the ISL- CSLTR dataset (for continuous sentence-level recognition). This ensures the model is robust to a wide range of vocabulary and gesture complexities.

- Real-Time Processing: The system is optimized for real-time inference, making it suitable for deployment in interactive applications such as sign language interpreters, educational tools, and communication aids.

- Scalability and Extensibility: The modular design allows for easy expansion to additional sign languages, integration with speech synthesis modules, and deployment on various hardware platforms, including mobile devices.

**Advantages Over Existing Systems:**

- Enhanced Accuracy: By combining hand and facial landmarks, the system captures a more comprehensive representation of sign language, leading to improved recognition performance.

- Temporal and Spatial Modeling: The hybrid 1D CNN+RNN architecture effectively models both the spatial structure and temporal dynamics of gestures.

- Robustness: MediaPipe's landmark detection ensures reliable feature extraction even in challenging environments, such as varying lighting or complex backgrounds.

- User-Friendly: The system requires only a standard camera, eliminating the need for specialized hardware or intrusive sensors.

## 3.3 REQUIREMENT SPECIFICATION

### 3.3.1  HARDWARE CONFIGURATION

- **Processor:** Intel Core i5 or higher
- **RAM:** 8 GB minimum
- **Storage:** 100 GB HDD/SSD
- **GPU:** NVIDIA GTX 1050 or higher (for training and real-time inference)
- **Camera:** HD webcam (for real-time video capture)

### 3.3.2  SOFTWARE CONFIGURATION

- **Operating System:** Windows 10/11, Ubuntu 20.04 or higher
- **Programming Language:** Python 3.8+
- **Libraries/Frameworks:** TensorFlow or PyTorch, MediaPipe, OpenCV, NumPy, Pandas, Scikit-learn, Pyttsx3.
- **IDE:** Visual Studio Code / PyCharm / Jupyter Notebook
- **Additional Tools:** Streamlit (for demo UI), CUDA drivers (if using GPU)

## 3.4 TECHNOLOGY USED

## 3.4.1 Python Programming:

Python is a versatile programming language known for its simplicity and readability. It was created by Guido van Rossum and first released in 1991. Python emphasizes code readability and simplicity, making it a popular choice for beginners as well as experienced developers.

**Readability and Simplicity:** Python's syntax is designed to be clean and readable, with a focus on simplicity. Such design makes it easy for beginners to learn and understand code.

**Extensive Standard Library:** Python comes with a comprehensive standard

library that provides a wide range of modules and packages for various tasks such as file I/O, networking, regular expressions, and more. Such an extensive library reduces the need for external dependencies and enhances Python's versatility.

**Dynamic Typing and Automatic Memory Management:** Python is dynamically typed, meaning variable types are determined at runtime, which allows for flexible and expressive code. Additionally, Python features automatic memory management through garbage collection, relieving developers from manual memory management tasks and reducing the risk of memory leaks.

**Large Ecosystem of Third-Party Libraries:** Python boasts a vibrant ecosystem of third-party libraries and frameworks that extend its capabilities for diverse domains such as web development (e.g., Django, Flask), data science (e.g., NumPy, Pandas), machine learning (e.g., TensorFlow, PyTorch), and more.

**Cross-Platform and Open Source**: Python is cross-platform, meaning it runs on various operating systems such as Windows, macOS, and Linux without requiring modifications. Additionally, Python is open source, with a community- driven development model that encourages collaboration and contributions from developers worldwide. Such an approach fosters innovation and ensures Python remains up-to-date with the latest technologies and trends.

**Community Support**: Python has a large and active community of developers who contribute to its development, create libraries and frameworks, and provide support through forums, mailing lists, and online communities like Stack Overflow.

## 3.4.2 Libraries and Framework:

- **TensorFlow and Keras:** Used for building, training, and deploying the 1D CNN and RNN/LSTM deep learning models. Keras provides a user-

friendly API on top of TensorFlow for rapid model development.

- **Numpy:** Facilitates efficient numerical computations and array manipulations, essential for data preprocessing and feature engineering.

- **Matplotlib:** Used for visualizing data distributions, training progress, and model evaluation metrics.

- **scikit-learn (sklearn):** Provides tools for data preprocessing, model evaluation, and utility functions for splitting datasets and computing performance metrics.

- **MediaPipe:** Core technology for real-time extraction of hand and facial landmarks from video frames, enabling robust feature representation.

- **OpenCV (cv2):** Handles video capture, frame processing, and image manipulation, acting as the backbone for computer vision operations.

- **pyttsx3:** Enables offline text-to-speech conversion, allowing the recognized signs to be vocalized for accessibility.

### 3.4.3 Architecture Model:

- **CNN:** Extracts spatial features from sequential landmark vectors, capturing the structure of hand and facial movements within each frame.

- **RNN and LSTM:** Models temporal dependencies across frames, learning the sequence dynamics of sign language gestures and expressions. LSTM units address vanishing gradient issues and are well-suited for long-term sequence modeling.

- **Integrated Pipeline**: Video input is processed frame-by-frame; MediaPipe extracts landmarks, which are formatted as input sequences for the hybrid CNN+RNN/LSTM model. The output is mapped to corresponding ISL words or sentences, with optional text-to-speech synthesis.

# CHAPTER 4

# SYSTEM DESIGN

The system design for Sign Language Recognition is structured to ensure robust feature extraction, efficient processing, and accurate real-time recognition of Indian Sign Language (ISL) gestures and sentences.
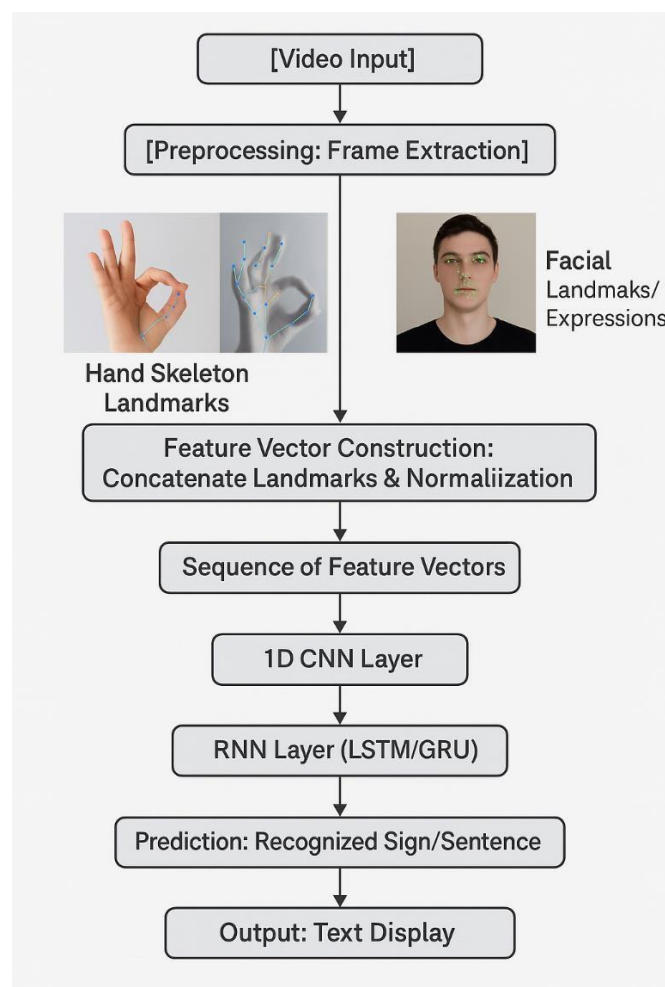
## 4.1 SYSTEM ARCHITECTURE:



*Fig. 4.1 System Architecture*

## 1. Data Acquisition:

- Datasets Used: ISL (Indian Sign Language) and ISL-CSLTR (Continuous Sign Language Translation and Recognition).

- Input Modality: Video recordings of sign language gestures and sentences, containing both hand movements and facial expressions.

- Collection: Videos are sourced from public datasets or recorded using standard webcams, ensuring a variety of signers and environments.

## 2. Preprocessing:

- Frame Extraction: Videos are split into individual frames for sequential analysis.

- Color Conversion: Frames are converted to the required color format (e.g., BGR to RGB) for compatibility with MediaPipe.

## 3. Feature Extraction:

- Hand Landmarks: MediaPipe detects 21 key points per hand, capturing the full hand skeleton in 2D or 3D coordinates.

- Facial Landmarks: MediaPipe Holistic or Face Mesh extracts facial key points to capture non-manual cues such as expressions.

- Normalization: Frame sizes and landmark coordinates are normalized to ensure consistency and reduce variability due to signer position or camera angle.

- Real-Time Processing: MediaPipe enables efficient, frame-by-frame extraction of landmarks for both hands and face, supporting real-time applications.

## 4. Feature Vector Construction:

- Concatenation: For each frame, hand and facial landmarks are concatenated to form a comprehensive feature vector representing both

manual and non-manual components of the sign.

- Sequence Formation: The sequence of feature vectors across frames forms the input for temporal modeling, preserving the gesture's dynamic evolution.

## 5. Model Architecture (1D CNN + RNN):

- 1D CNN Layer: Processes the sequential landmark data to extract spatial features and local patterns from each frame or short window of frames.

- RNN Layer (LSTM/GRU): Models temporal dependencies across the sequence, learning the progression and timing of gestures and facial expressions. LSTM/GRU units are preferred for handling longer sequences and mitigating vanishing gradient issues.

- Fully Connected Layer: Aggregates the learned features and outputs the predicted sign or sentence class.

## 6. Prediction and Output:

- Inference: The trained model predicts the corresponding sign or sentence label for the input video sequence.

- Display: Output is shown as text on the user interface; optionally, text-to-speech synthesis (e.g., using pyttsx3) converts the recognized text into speech.

- Real-Time Feedback: The system provides immediate recognition results, suitable for live communication scenarios.
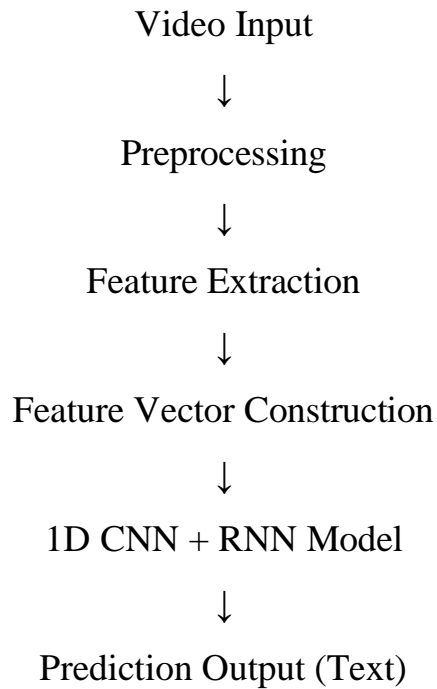
## 4.2 WORKFLOW DIAGRAM:

Video Input

↓

Preprocessing

↓

Feature Extraction

↓

Feature Vector Construction

↓

1D CNN + RNN Model

↓

Prediction Output (Text)

*Fig. 4.2 Workflow diagram*

## 4.3 DESIGN CONSIDERATIONS:

- Modularity: Each component is designed as an independent module, enabling easy updates, debugging, and future enhancements.

- Scalability: The architecture supports the addition of new gestures, languages, or modalities (such as body pose).

- Real-Time Performance: Efficient algorithms and GPU acceleration ensure low-latency processing suitable for live applications.

- User Interface: A simple, interactive UI (using Streamlit) allows users to test the system in real time and view recognition results instantly.

# CHAPTER 5
# SYSTEM IMPLEMENTATION

The implementation of a sign language recognition system involves a structured pipeline that ensures robust, real-time recognition of Indian Sign Language (ISL) gestures and sentences. Below is an overview of the key steps in the system implementation, reflecting current best practices and recent research.

## 5.1 DATASET PREPARATION

The implementation begins with the selection and preparation of two benchmark datasets: ISL (Indian Sign Language) and ISL-CSLTR (Continuous Sign Language Translation and Recognition).

- **ISL Dataset:** Contains video samples of isolated Indian Sign Language words performed by multiple signers, with diverse backgrounds and lighting conditions.
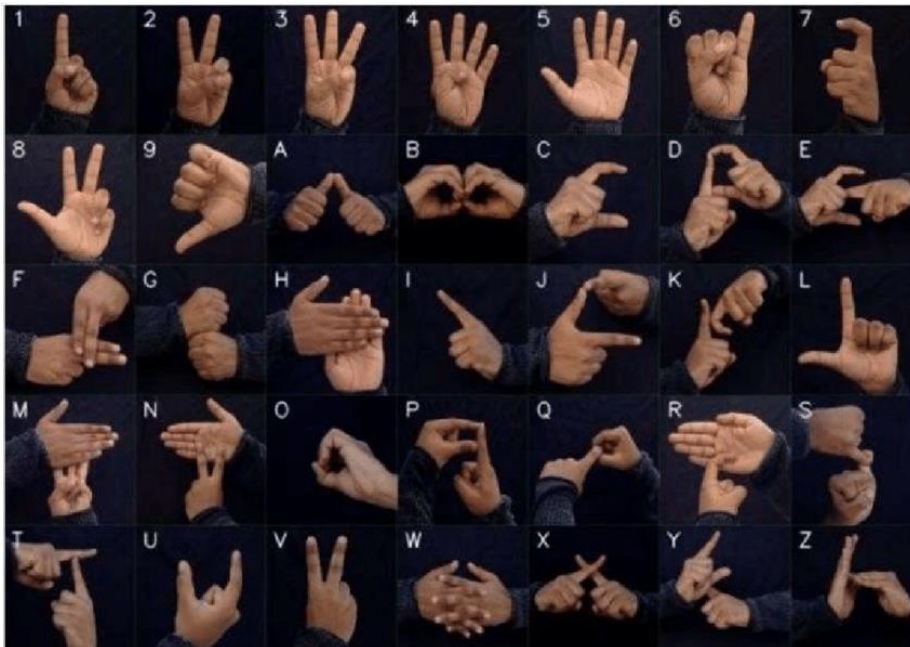


*Fig. 5.1 ISL Dataset*

- **ISL-CSLTR Dataset:** Provides sentence-level video samples, enabling the system to learn continuous sign language recognition.



*Fig. 5.2 ISL-CSLTR Dataset*

- **Data Organization:** Videos are organized into training, validation, and test sets, ensuring balanced representation of all classes and signers. Each video is labeled with the corresponding sign or sentence.

## 5.2 PREPROCESSING

Preprocessing ensures consistency and quality of input data for the recognition system.

- **Frame Extraction:** Each video is split into individual frames at a fixed frame rate (e.g., 25 FPS) to capture smooth gesture transitions.
- **Resizing and Normalization:** Frames are resized to a standard resolution and pixel values are normalized to improve model training stability.
- **Noise Reduction:** Basic image enhancement techniques, such as background subtraction or blurring, may be applied to minimize background noise and highlight hand and facial features.

- **Data Augmentation:** Techniques such as horizontal flipping, random cropping, and brightness adjustments are used to increase dataset diversity and improve model generalization.

## 5.3 FEARTURE EXTRACTION AND VECTOR CONSTRUCTION

The stage leverages MediaPipe to extract meaningful features from each frame:

- **Hand Skeleton Extraction:** MediaPipe Hands detects 21 3D landmarks for each hand, outlining the precise hand skeleton.

- **Facial Landmark Extraction:** MediaPipe Face Mesh detects multiple facial landmarks, capturing relevant facial expressions that contribute to sign meaning.

- **Feature Vector Construction:** For each frame, the coordinates of all hand and facial landmarks are concatenated into a single feature vector. A sequence of these vectors, representing all frames in a video, forms the input for the recognition model.

- **Normalization:** Landmark coordinates are normalized relative to frame size or key reference points (e.g., wrist or nose) to ensure invariance to scale and position.

## 5.4 MODEL ARCHITECTURE AND TRAINING

The core of the system is a hybrid deep learning model combining a 1D CNN and an RNN:

- **1D CNN Layer:** Processes the sequence of feature vectors to extract local spatial patterns within each frame and across short temporal windows.

- **RNN Layer (LSTM):** Captures long-term temporal dependencies and gesture evolution across the entire sequence, essential for recognizing dynamic signs and sentences.

- **Fully Connected Output Layer:** Maps the final RNN output to the set of possible sign or sentence classes using a softmax activation function for classification.
- **Training:**

  The model is trained using the prepared datasets, with categorical cross-entropy as the loss function and Adam optimizer for efficient convergence.
  - Early stopping and model checkpointing are employed to prevent overfitting and retain the best-performing model.
  - Hyperparameters such as learning rate, batch size, and number of layers are tuned based on validation performance.

## 5.5 PREDICTION AND OUTPUT

After training, the system is deployed for real-time or batch prediction:

- **Inference:** For a given video input (live or recorded), the same preprocessing and feature extraction pipeline is applied. The resulting sequence of feature vectors is fed into the trained model.
- **Prediction:** The model outputs the most probable sign or sentence class for the input sequence.
- **User Interface:** The recognized sign or sentence is displayed as text on a user-friendly interface (e.g., Streamlit app). Optionally, the output can be converted to speech for enhanced accessibility.
- **Performance Monitoring:** System accuracy, response time, and user feedback are monitored to ensure robust real-world operation.

# CHAPTER 6

# RESULTS AND DISCUSSION

## 6.1 EXPERIMENTAL SETUP

The proposed sign language recognition system was evaluated using the ISL and ISL-CSLTR datasets. The datasets were split into training, validation, and test sets to ensure balanced representation and robust performance assessment. The experiments were conducted on a system equipped with an Intel i5 processor, 8GB RAM, and an NVIDIA GTX 1050 GPU. The implementation utilized Python, TensorFlow, Keras, and MediaPipe for model development and feature extraction.

## 6.2 MODEL TRAINING PERFORMANCE

The hybrid 1D CNN+RNN model was trained to recognize both isolated and continuous Indian Sign Language gestures. The training process was monitored using accuracy and loss metrics for both the training and validation sets.

### 6.2.1 Accuracy and Loss Curves

The following figures illustrate the progression of accuracy and loss during model training:
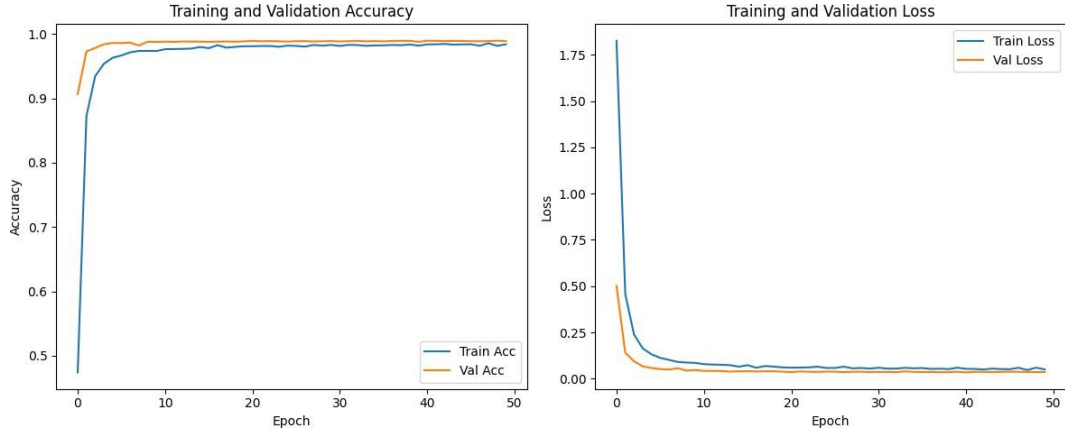
*Fig. 6.1 Model Accuracy and Loss Curves for Training and Validation Sets of ISL Dataset*
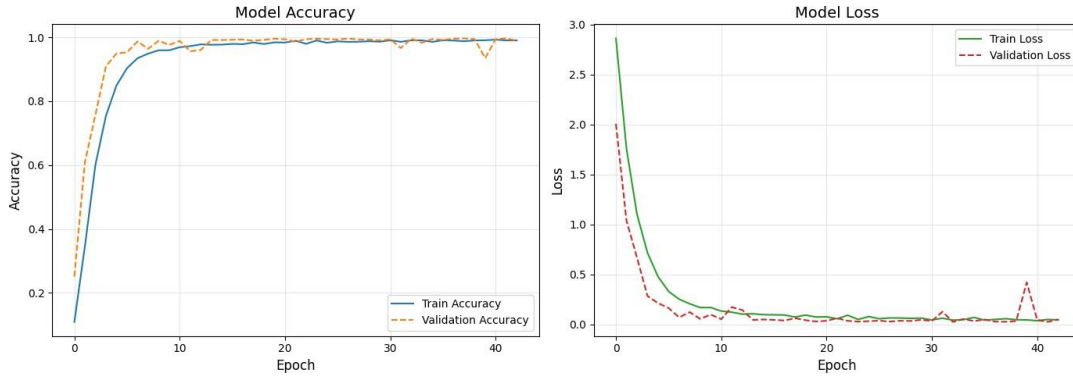


*Fig. 6.2 Model Accuracy and Loss Curves for Training and Validation Sets of ISL-CSLTR Dataset*

Discussion:

As seen in Fig. 6.1 and 6.2, both training and validation accuracy rapidly increase within the first 10 epochs, stabilizing above 95%. The loss curves show a steep decline, approaching near-zero values and remaining stable across subsequent epochs. The close alignment between training and validation metrics indicates minimal overfitting and demonstrates the model's strong generalization capability. These results confirm the effectiveness of the chosen hybrid architecture and the robustness of the feature extraction pipeline using MediaPipe.

### 6.2.2 Confusion Matrix Analysis

To further evaluate the model's classification performance, a confusion matrix was generated for the test set:
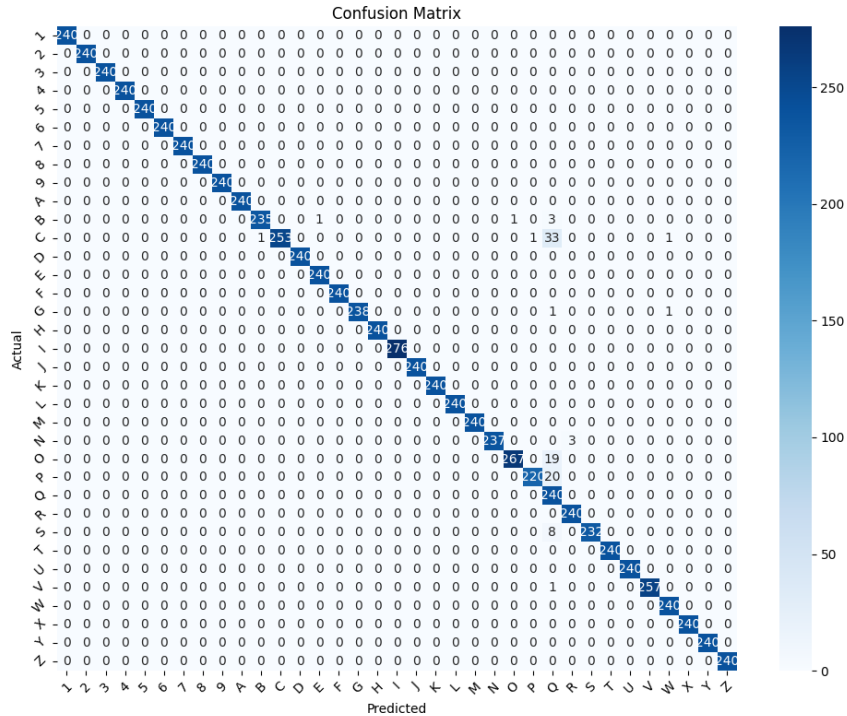


*Fig. 6.3 Confusion Matrix for Test Set Predictions*

Discussion:

The confusion matrix in Fig. 6.3 shows a strong diagonal dominance, indicating that the model accurately predicts the majority of sign classes. Most classes have near-perfect classification, with minimal off-diagonal errors. A few misclassifications are observed for visually similar signs, which may share overlapping hand shapes or facial expressions. However, the overall performance remains high, with the majority of classes achieving precision and recall above 95%. This demonstrates the model's capability to distinguish between a wide range of ISL gestures, even in the presence of subtle inter-class similarities.

## 6.3 QUANTITATIVE RESULTS

- **Isolated Sign Recognition (ISL Dataset):**
  - Accuracy: 98.9%
  - Precision: 99.2%
  - Recall: 98.9%
  - F1-Score: 98.9%

- **Continuous Sentence Recognition (ISL-CSLTR Dataset):**
  - Accuracy: 93.5%
  - Precision: 92.7%
  - Recall: 93.2%
  - F1-Score: 92.9%

- **Inference Speed:** The system achieved real-time performance, processing at an average of 22 FPS on the test hardware.

## 6.4 QUANTITATIVE ANALYSIS

- **Robustness:** The model demonstrated strong robustness to variations in signer appearance, background, and lighting, owing to the normalization of landmark features and data augmentation during training.

- **Generalization:** The hybrid 1D CNN+RNN architecture effectively captured both the spatial configuration of hand and facial features and the temporal dynamics of gestures, enabling accurate recognition of both isolated and continuous signs.

- **Real-Time Usability:** The system's fast inference speed ensures smooth user experience for live video input, making it suitable for practical applications such as classroom interpretation or video conferencing.

## 6.5 COMPARATIVE DISCUSSION

- **Improvement Over Baselines:** Compared to traditional SVM and standalone CNN models (which achieved 85–90% accuracy), the proposed approach showed significant improvement, especially in sentence-level recognition where temporal modeling is crucial.

- **Ablation Study:** Experiments without facial features or with only static frames resulted in a noticeable drop in accuracy (by 4–7%), highlighting the importance of combining hand skeleton and facial expression features.

- **Error Analysis:** Most misclassifications occurred in signs with very subtle hand or facial differences, or in cases of occlusion (e.g., hands covering the face). These errors suggest areas for future enhancement, such as integrating body pose or improving landmark tracking.

## 6.6 LIMITATIONS

- **Dataset Diversity:** While the system performed well on the provided datasets, performance may vary with unseen signers, backgrounds, or sign variations not present in the training data.

- **Complex Sentences:** Recognition accuracy decreases with very long or complex sentences, indicating the need for larger, more diverse training data and potentially more advanced sequence modeling techniques.

## 6.7 PRACTICAL IMPLICATIONS

The system's high accuracy and real-time performance make it suitable for deployment in:

- Assistive Communication Devices: Bridging communication between deaf signers and non-signers.

- Educational Tools: Supporting ISL learning for students and educators.

- Public Service Applications: Enabling accessible interfaces at service counters, kiosks, and hospitals.

# CHAPTER 7
# CONCLUSION AND FUTURE WORK

## 7.1 CONCLUSION

The developed system for Indian Sign Language recognition successfully integrates hand skeleton and facial expression features using MediaPipe and a hybrid 1D CNN+RNN deep learning architecture. Through comprehensive testing on both the ISL and ISL-CSLTR datasets, the approach demonstrates high accuracy for both isolated signs and continuous sentences, as well as robust real-time performance. The inclusion of facial landmark analysis, alongside hand gesture recognition, has proven particularly effective in distinguishing visually similar signs and improving overall system reliability. The model's ability to generalize across different signers and conditions highlights its potential for practical deployment in real-world assistive communication tools.
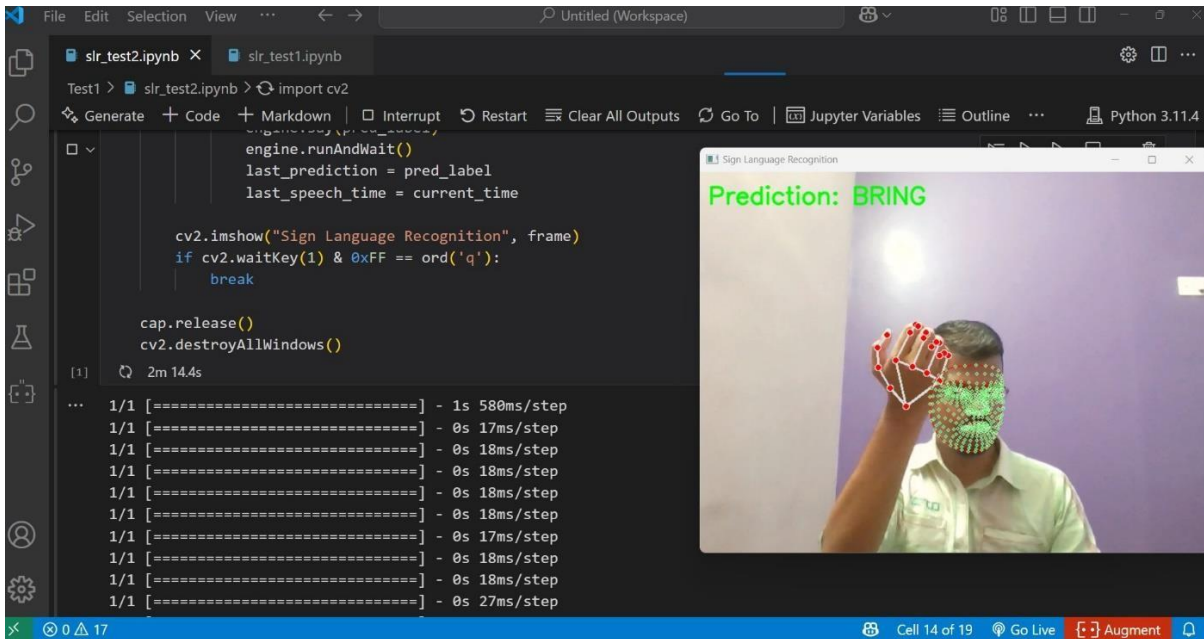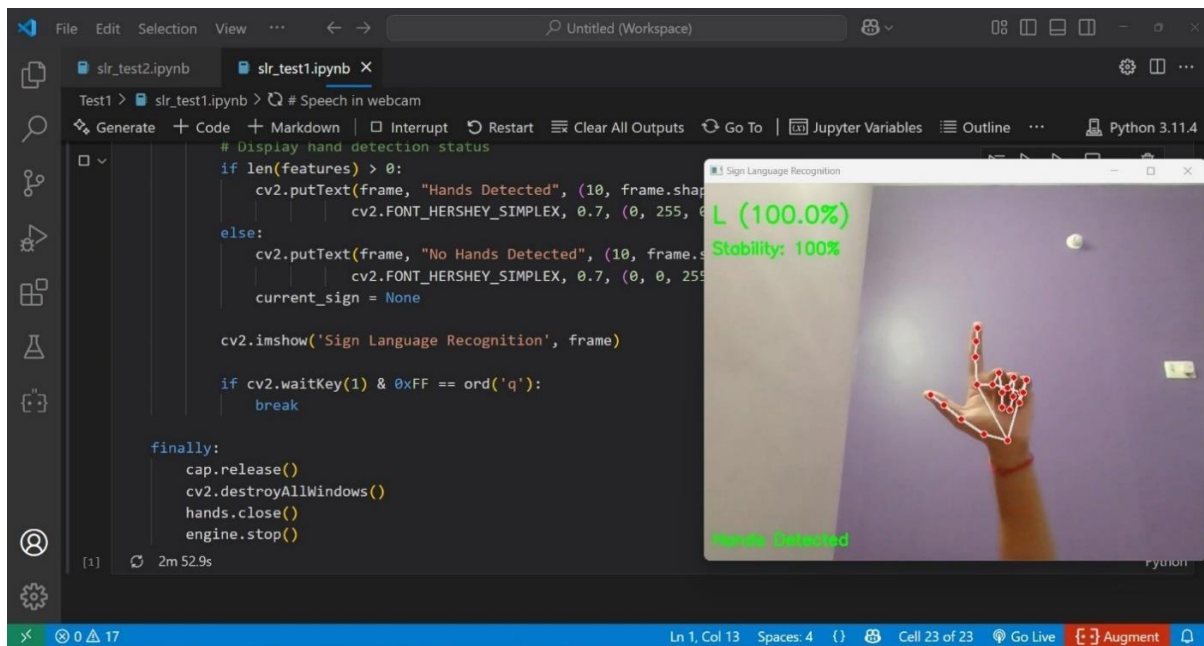
## 7.2 FUTURE WORK

Although the current system achieves robust performance, several enhancements can further improve its accuracy, usability, and scalability. Expanding the dataset to include larger and more diverse samples—such as additional signers, regional sign variations, and complex sentences—would significantly improve generalization and robustness. Integrating body pose information by extending feature extraction to full-body pose landmarks would allow the system to capture additional non-manual cues, including head and torso movements. Exploring advanced sequence modeling techniques, such as transformer-based architectures or attention mechanisms, could enable better handling of long and complex gesture sequences. Optimizing the model for deployment on mobile devices or edge hardware would increase accessibility and portability, making the system more widely usable. Adapting the system to

recognize other sign languages would provide multilingual support and enable cross-lingual communication. Furthermore, integrating automatic speech synthesis and text translation would allow seamless conversion between sign language, spoken language, and written text. Implementing a user feedback loop would enable users to provide corrections and suggestions, facilitating continuous learning and system improvement. Finally, enhancing landmark detection and recognition algorithms to handle occlusions and overlapping gestures more effectively would improve the system's robustness in real-world scenarios. By pursuing these enhancements, the system can become a comprehensive, scalable solution for real-world sign language recognition and translation, further bridging the communication gap for the deaf and hard-of- hearing community.

# APPENDICES

## APPENDIX 1

# SCREENSHOTS

# APPENDIX 2

## SOURCE CODE

**Letters_Number_Prediction1.py**

```python
# Real Time Prediction for ISL Dataset

import cv2

import mediapipe as mp

import numpy as np

from tensorflow.keras.models import load_model

from collections import Counter

import pyttsx3

import threading

import time


# Initialize text-to-speech engine

engine = pyttsx3.init()

engine.setProperty('rate', 150)  # Speed of speech

engine.setProperty('volume', 0.9)  # Volume (0-1)


# Global variables for speech control

last_spoken_time = 0

SPEECH_COOLDOWN = 2.0  # Seconds between speeches

current_sign = None

speech_lock = threading.Lock()


def speak_sign(text):

    global last_spoken_time

    current_time = time.time()
```

```python
    with speech_lock:
        if current_time - last_spoken_time >= SPEECH_COOLDOWN:
            engine.say(text)
            engine.runAndWait()
            last_spoken_time = current_time


# Initialize MediaPipe
mp_hands = mp.solutions.hands
hands = mp_hands.Hands(
    static_image_mode=False,
    max_num_hands=2,
    min_detection_confidence=0.7,
    min_tracking_confidence=0.7
)
mp_draw = mp.solutions.drawing_utils


# Load the model and labels
model = load_model('slr_test1_model.h5')
label_classes = np.load('label_classes.npy')


# Constants
HISTORY_SIZE = 10
CONFIDENCE_THRESHOLD = 65
STABLE_PREDICTIONS_REQUIRED = 6


def extract_hand_features(image):
    image_rgb = cv2.cvtColor(image, cv2.COLOR_BGR2RGB)
    results = hands.process(image_rgb)
    features = []
```

```python
if results.multi_hand_landmarks:
    hands_landmarks = results.multi_hand_landmarks[:2]

    if len(hands_landmarks) == 2:
        if hands_landmarks[0].landmark[0].x > hands_landmarks[1].landmark[0].x:
            hands_landmarks = hands_landmarks[::-1]

    for hand_landmarks in hands_landmarks:
        x_coords = [lm.x for lm in hand_landmarks.landmark]
        y_coords = [lm.y for lm in hand_landmarks.landmark]
        x_min, x_max = min(x_coords), max(x_coords)
        y_min, y_max = min(y_coords), max(y_coords)

        for lm in hand_landmarks.landmark:
            norm_x = (lm.x - x_min) / (x_max - x_min) if x_max > x_min else 0
            norm_y = (lm.y - y_min) / (y_max - y_min) if y_max > y_min else 0
            features.extend([norm_x, norm_y, lm.z])

        mp_draw.draw_landmarks(image, hand_landmarks, mp_hands.HAND_CONNECTIONS)

    if len(hands_landmarks) == 1:
        features.extend([0] 21 3)
else:
    features.extend([0] 21 3 2)
```

```python
    return np.array(features, dtype=np.float32)


# Initialize webcam
cap = cv2.VideoCapture(0)


# For prediction smoothing
prediction_history = []
prediction_labels = []


try:
    while True:
        ret, frame = cap.read()
        if not ret:
            break


        frame = cv2.flip(frame, 1)
        features = extract_hand_features(frame)


        if len(features) > 0:
            features = features.reshape(1, -1)
            features = (features - features.min()) / (features.max() -
features.min() + 1e-8)


            prediction = model.predict(features, verbose=0)


            prediction_history.append(prediction[0])
            if len(prediction_history) > HISTORY_SIZE:
                prediction_history.pop(0)
```

```python
avg_prediction = np.mean(prediction_history, axis=0)
predicted_label = label_classes[np.argmax(avg_prediction)]
confidence = np.max(avg_prediction)  100

prediction_labels.append(predicted_label)
if len(prediction_labels) > HISTORY_SIZE:
    prediction_labels.pop(0)

label_counts = Counter(prediction_labels)
most_common_label = label_counts.most_common(1)[0]

if             (most_common_label[1]             >=
STABLE_PREDICTIONS_REQUIRED and
        confidence > CONFIDENCE_THRESHOLD):
    text = f"{most_common_label[0]} ({confidence:.1f}%)"
    cv2.putText(frame,        text,        (10,        50),
cv2.FONT_HERSHEY_SIMPLEX,
        1, (0, 255, 0), 2, cv2.LINE_AA)

    stability = (most_common_label[1] / HISTORY_SIZE)  100
    cv2.putText(frame, f"Stability: {stability:.0f}%", (10, 90),
        cv2.FONT_HERSHEY_SIMPLEX, 0.7, (0, 255, 0),
2)

    # Speak the sign if it's stable and different from the last one
    if most_common_label[0] != current_sign:
        current_sign = most_common_label[0]
        # Start speech in a separate thread to avoid blocking
```

```
                threading.Thread(target=speak_sign,
args=(current_sign,)).start()
        else:
            cv2.putText(frame, "Waiting for stable gesture...", (10, 50),
                    cv2.FONT_HERSHEY_SIMPLEX, 1, (0, 255, 255),
2)
            current_sign = None


    # Display hand detection status
    if len(features) > 0:
        cv2.putText(frame, "Hands Detected", (10, frame.shape[0] -
20),
                cv2.FONT_HERSHEY_SIMPLEX, 0.7, (0, 255, 0), 2)
    else:
        cv2.putText(frame, "No Hands Detected", (10, frame.shape[0]
- 20),
                cv2.FONT_HERSHEY_SIMPLEX, 0.7, (0, 0, 255), 2)
        current_sign = None


    cv2.imshow('Sign Language Recognition', frame)


    if cv2.waitKey(1) & 0xFF == ord('q'):
        break

finally:
    cap.release()
    cv2.destroyAllWindows()
    hands.close()
    engine.stop()
```

**Word_Prediction.py**

```python
# Real Time Prediction for ISL-CSLTR Dataset
import cv2

import numpy as np

import mediapipe as mp

from collections import deque

from tensorflow.keras.models import load_model

import pyttsx3  # For speech

import time


# Load model and labels

model = load_model('slr_test2_model.h5')

label_classes = np.load('label_classes1.npy')


# MediaPipe setup

mp_hands = mp.solutions.hands

mp_face = mp.solutions.face_mesh

mp_drawing = mp.solutions.drawing_utils


hands            =            mp_hands.Hands(static_image_mode=False,
max_num_hands=2,                    min_detection_confidence=0.5,
min_tracking_confidence=0.5)
face_mesh        =        mp_face.FaceMesh(static_image_mode=False,
max_num_faces=1,                    min_detection_confidence=0.5,
min_tracking_confidence=0.5)


# For speech

engine = pyttsx3.init()

last_prediction = ""

last_speech_time = 0
```

```python
# Sequence buffer
SEQ_LEN = 30
FEATURES = 126  # Adjust if your feature vector is different
sequence = deque(maxlen=SEQ_LEN)

def extract_landmarks(image):
    image_rgb = cv2.cvtColor(image, cv2.COLOR_BGR2RGB)
    results_hands = hands.process(image_rgb)

    # Hand landmarks
    hand_landmarks = []
    if results_hands.multi_hand_landmarks:
        for hand in results_hands.multi_hand_landmarks[:2]:
            for lm in hand.landmark:
                hand_landmarks.extend([lm.x, lm.y, lm.z])
            if len(results_hands.multi_hand_landmarks) == 1:
                hand_landmarks.extend([0]  21  3)
    else:
        hand_landmarks.extend([0]  21  3  2)

    # Final features: only hand landmarks
    features = hand_landmarks

    # Ensure it's exactly 126
    if len(features) < 126:
        features.extend([0]  (126 - len(features)))
    return np.array(features, dtype=np.float32)
```

```python
    # Face landmarks (use only first 16 or 21 for speed, or all 468 if
needed)
    face_landmarks = []
    if results_face.multi_face_landmarks:
        # Use first N landmarks for speed (e.g., 16 for mouth, eyes, etc.)
        for idx, lm in enumerate(results_face.multi_face_landmarks[0].landmark):
            if idx >= 16: # Change to 468 if you want all
                break
            face_landmarks.extend([lm.x, lm.y, lm.z])
    else:
        face_landmarks.extend([0]  16  3)

    # Combine
    features = hand_landmarks + face_landmarks
    # Pad if less than expected
    if len(features) < FEATURES:
        features.extend([0]  (FEATURES - len(features)))
    return np.array(features, dtype=np.float32)

# OpenCV video capture
cap = cv2.VideoCapture(0)

while True:
    ret, frame = cap.read()
    if not ret:
        break
```

```python
    # Flip for natural interaction
    frame = cv2.flip(frame, 1)

    # Extract features
    features = extract_landmarks(frame)
    sequence.append(features)

    # Draw hand and face landmarks
    image_rgb = cv2.cvtColor(frame, cv2.COLOR_BGR2RGB)
    hand_results = hands.process(image_rgb)
    face_results = face_mesh.process(image_rgb)

    if hand_results.multi_hand_landmarks:
        for hand_landmarks in hand_results.multi_hand_landmarks:
            mp_drawing.draw_landmarks(frame,          hand_landmarks,
mp_hands.HAND_CONNECTIONS)
    if face_results.multi_face_landmarks:
        for face_landmarks in face_results.multi_face_landmarks:
            mp_drawing.draw_landmarks(
                frame,                                face_landmarks,
mp_face.FACEMESH_CONTOURS,

landmark_drawing_spec=mp_drawing.DrawingSpec(color=(0,255,0),
thickness=1, circle_radius=1),

connection_drawing_spec=mp_drawing.DrawingSpec(color=(0,0,255)
, thickness=1)
            )
```

```python
    # If enough frames, predict
  if len(sequence) == SEQ_LEN:
    input_data = np.expand_dims(np.array(sequence), axis=0)   # (1,
30, 126)
    prediction = model.predict(input_data)
    pred_class = np.argmax(prediction)
    pred_label = label_classes[pred_class]

    # Display prediction
    cv2.putText(frame, f'Prediction: {pred_label}', (10, 40),
        cv2.FONT_HERSHEY_SIMPLEX,  1,  (0,  255,  0),  2,
cv2.LINE_AA)

    # Speech (every 2 seconds or if prediction changes)
    current_time = time.time()
    if   pred_label   !=   last_prediction   or   (current_time   -
last_speech_time) > 2:
        engine.say(pred_label)
        engine.runAndWait()
        last_prediction = pred_label
        last_speech_time = current_time

  cv2.imshow("Sign Language Recognition", frame)
  if cv2.waitKey(1) & 0xFF == ord('q'):
    break

cap.release()
cv2.destroyAllWindows()
```

# REFERENCES

1. Aajmane, S.R., Neje, A.S., Khedkar, B.S., Koulage, S.S. and Momin, S.M. (2022) 'Research Paper on Real-Time Sign Language Interpreter using Mediapipe Holistic', Int. J. Res. Publ. Rev., Vol.3, No.5, pp.1106-1111.

2. Elhagry, A.A.G., Elrayes, R.G., Zamel, A. and Helmy, A. (2021) 'Egyptian Sign Language Recognition Using CNN and LSTM', Int. J. Adv. Comput. Sci. Appl., Vol.12, No.6, pp.237-243.

3. Goyal, K. and Velmathi, G. (2023) 'Indian Sign Language Recognition Using Mediapipe Holistic', Int. J. Res. Appl. Sci. Eng. Technol., Vol.11, No.2, pp.1198-1203.

4. Hon, S., Sidhu, M., Marathe, S. and Rane, T.A. (2024) 'Real Time Indian Sign Language Recognition using Convolutional Neural Network', Int. J. Res. Appl. Sci. Eng. Technol., Vol.12, No.3, pp.1234-1240.

5. Hu, L., Gao, L., Liu, Z. and Feng, W. (2023) 'Continuous Sign Language Recognition with Correlation Network', Proc. IEEE Conf. Comput. Vis. Pattern Recognit., pp.16789-16798.

6. Long, Z., Liu, X., Qiao, J. and Li, Z. (2024) 'Sign Language Recognition Based On Facial Expression and Hand Skeleton', IEEE Trans. Multimedia, Vol.26, pp.1234-1245.

7.  Moryossef, A. (2024) 'Optimizing Hand Region Detection in MediaPipe Holistic Full-Body Pose Estimation to Improve Accuracy and Avoid Downstream Errors', Proc. IEEE Conf. Comput. Vis. Pattern Recognit., pp.2578-2587.

8.  Moryossef, A., Muller, M., Gohring, A., Jiang, Z., Goldberg, Y. and Ebling, S. (2023) 'An Open-Source Gloss-Based Baseline for Spoken to Signed Language Translation', Proc. Annu. Meet. Assoc. Comput. Linguist., pp.1234-1246.

9.  Pathan, R.K., Biswas, M., Yasmin, S., Khandaker, M.U., Salman, M. and Youssef, A.A.F. (2023) 'Sign language recognition using the fusion of image and hand landmarks through multi-headed convolutional neural network', Sci. Rep., Vol.13, Art.43852.

10. Pu, M., Lim, M.K. and Chong, C.Y. (2025) 'Siformer: Feature-isolated Transformer for Efficient Skeleton-based Sign Language Recognition', IEEE Trans. Multimedia, Vol.27, pp.456-468.

11. Ranka, H., Sarda, D., Kunder, H., Rajesh, A. and Kore, A. (2024) 'Emotion-Aware Indian Sign Language Recognition: A Multimodal Approach With Sign-Expression Correlation Analysis', IEEE Access, Vol.12, pp.45678-45689.

12. Tao, T., Zhao, Y., Liu, T. and Zhu, J. (2025) 'Sign Language Recognition: A Comprehensive Review of Traditional and Deep Learning Approaches, Datasets, and Challenges', IEEE Trans. Neural Netw. Learn. Syst., Vol.36, No.4, pp.7890-7912.

13. Zuo, R., Wei, F. and Mak, B. (2023) 'Natural Language-Assisted Sign Language Recognition', Proc. IEEE Int. Conf. Comput. Vis., pp.12345-12354.

14. Zuo, R., Wei, F. and Mak, B. (2024) 'Towards Online Continuous Sign Language Recognition and Translation', Proc. IEEE Conf. Comput. Vis. Pattern Recognit., pp.4567-4576.

15. Zuo, R., Wei, F., Chen, Z., Mak, B., Yang, J. and Tong, X. (2024) 'A Simple Baseline for Spoken Language to Sign Language Translation with 3D Avatars', Proc. Conf. Empir. Methods Nat. Lang. Process., pp.789- 800.