# AI Tool for Indian Sign Language Generator from Audio-Visual Content in English/Hindi and Vice-Versa

**Aishwarya K A, Arpitha M, Bindu P, and G Poornima**

*{aishwaryaka2104, arpitham1893, bindupthumma, gpoornimabdvt}@gmail.com*

*Department of Computer Science and Engineering, PES Institute of Technology and Management, Shivamogga, India*

*Abstract*—People with hearing or speech impairments often face communication barriers in their day-to-day lives.The need for accessible communication tools for individuals with hearing or speech impairments continues to grow as digital communication and technology-driven interaction become central to everyday life. Indian Sign Language (ISL), being a primary mode of communication for many, requires automated systems that can interpret gestures accurately and efficiently. Motivated by this need, our project presents an AI-driven ISL Recognition and Translation System that integrates Computer Vision, Deep Learning, and Natural Language Processing (NLP) within a unified web-based platform. The system employs MediaPipe to extract reliable hand landmarks and uses a Convolutional Neural Network (CNN) to learn and classify spatial gesture patterns from webcam streams or uploaded videos. For text-to-sign translation, the framework applies NLP techniques—including tokenization, lemmatization, and tense detection—to generate meaningful ISL animations. When a direct gesture animation is not available, the system automatically breaks words into alphabet-level signs to ensure complete translation. The preprocessing pipeline includes region-of-interest extraction and normalization to maintain consistent input quality. Implemented using a Django web framework with an intuitive user interface, the system supports educational, assistive, and accessibility-focused applications. Overall, this project delivers a reliable, user-centered solution that enhances communication and promotes digital inclusivity for the hearing and speech-impaired community

Keywords: Indian Sign Language (ISL), Gesture Recognition, Convolutional Neural Network (CNN), MediaPipe, Natural Language Processing (NLP), Assistive Technology.

## I. INTRODUCTION

The growth of Artificial Intelligence (AI) and Computer Vision has transformed how gesture-based communication systems are designed, interpreted, and deployed. With the increasing reliance on digital communication, the need for accessible tools that assist individuals with hearing or speech impairments has become more critical than ever. Indian Sign Language (ISL), the primary mode of communication for many in the hearing-impaired community, requires sophisticated automated systems capable of interpreting hand gestures accurately in real time. Traditional approaches relied heavily on hand-crafted features and simple image-processing methods, which often struggled with variations in lighting, hand orientation, and background noise [1]. However, advancements in deep

learning have significantly improved recognition capabilities, making gesture interpretation more reliable and scalable in real-world settings.

Earlier ISL recognition systems primarily used basic image-processing techniques and rule-based models that produced low accuracy and were limited to controlled environments. With the introduction of powerful learning models, particularly Convolutional Neural Networks (CNNs), the field experienced remarkable progress as these models learned spatial features directly from gesture frames without requiring handcrafted rules [4]. Further improvements were made through the integration of Natural Language Processing (NLP), enabling the conversion of natural text or speech into ISL animations using linguistic techniques such as tokenization, lemmatization, and tense identification [2]. Recent models also explore hybrid architectures combining CNNs with temporal networks like LSTM or GRU to handle continuous sign sequences, improving recognition across dynamic gestures where motion consistency is crucial [5].

Despite these advancements, several challenges persist in sign language recognition. Variations in users' hand shapes, camera placement, background clutter, and inconsistent illumination significantly affect model performance. Moreover, most existing systems are limited either to static signs or controlled datasets, making them less effective in everyday real-world environments. Many approaches also focus solely on recognition and do not provide a complete translation mechanism that supports text-to-sign animations or fallback strategies when specific gestures are unavailable. These gaps highlight the need for an integrated framework that combines accuracy, flexibility, and user accessibility [6].

To address these limitations, this work proposes an AI-driven Indian Sign Language Recognition and Translation System capable of analyzing both static and dynamic gestures in real time. The system leverages MediaPipe for robust hand landmark detection and employs a CNN model to learn spatial gesture patterns from webcam and uploaded video inputs. Additionally, a text-to-sign module powered by NLP generates meaningful ISL animations, decomposing words into alphabet-level signs when direct gestures are missing [12]. The solution is deployed through a Django-based web

interface, making it practical for real-world use in educational, assistive, and accessibility-centered applications. Beyond improving recognition accuracy, this project emphasizes usability and scalability, contributing toward a more inclusive digital environment where technology bridges communication gaps instead of widening them.

## II. LITERATURE REVIEW

The paper titled "ISL Words Recognition Using Hybrid Transfer Learning and RNN" was authored by Naman Bansal and Abhilasha Jain [1]. They focused on developing a system for recognizing Indian Sign Language (ISL) words using both static and dynamic gesture data. The dataset used in their research consisted of image and video samples representing different ISL words. These samples were preprocessed through resizing frames, removing noise, and normalizing pixel values to maintain uniform input dimensions. In addition, video frames were converted into sequential data to facilitate temporal learning and improve feature extraction.

The main algorithm used in this work was a hybrid model combining Transfer Learning with a Recurrent Neural Network (RNN) classifier. Pre-trained CNN architectures such as VGG16 and InceptionV3 were employed for spatial feature extraction, while the RNN component was responsible for capturing the temporal sequence patterns present in gesture videos. This combination effectively utilized both spatial and temporal characteristics for accurate ISL word recognition.

The model achieved a performance accuracy of around 97–98%, demonstrating excellent generalization capability and lower error rates compared to conventional CNN-only models. This indicates that the system can efficiently recognize ISL words with high precision and reliability.

The paper titled "Harnessing AI to Generate Indian Sign Language from Natural Speech" was authored by Alisha Kulkarni et al. [2]. They focused on developing a system that translates spoken English into Indian Sign Language (ISL) gestures. The dataset used in their research consisted of custom-created ISL gesture samples corresponding to commonly used English words and phrases. The audio inputs were first processed through standard speech-to-text APIs to convert speech into textual data. This text was then tokenized and mapped to the appropriate ISL gestures. The gesture visuals—comprising video frames or image sequences—were resized and normalized to ensure uniform representation. Additional preprocessing steps included background noise reduction and segmentation to enhance gesture clarity and recognition accuracy.

The main algorithm used in this work followed a three-stage pipeline integrating speech recognition, text interpretation, and gesture rendering. The speech recognition module captured real-time voice input and converted it into text. The textual data was then semantically analyzed using Natural Language Processing (NLP) techniques to determine the intended meaning. Finally, a gesture visualization module generated the corresponding ISL signs using pre-recorded video snippets or 3D gesture animations. The entire system was implemented in Python with the support of machine learning frameworks to ensure high accuracy and efficiency.

The model achieved a performance accuracy of approximately 94–96%, demonstrating strong reliability in mapping English speech to ISL gestures. The system exhibited real-time responsiveness and low latency, proving effective for assistive communication, accessibility enhancement, and educational applications for the hearing-impaired community.

The paper titled "AI-Based Real-Time Indian Sign Language Recognition System Using Deep Learning" was authored by Kujani T. and Dhilip Kumar V. [3]. They focused on developing a deep learning-based system for translating Indian Sign Language (ISL) gestures to assist individuals with verbal impairment. The dataset used in their research comprised 35 ISL hand gestures representing alphabets and commonly used words. Thousands of images were captured under different lighting and background conditions to ensure variety and robustness. The preprocessing steps included image resizing, contrast enhancement, background noise removal, normalization, and grayscale conversion to improve gesture clarity and reduce computational complexity.

The main algorithm used in this work involved multiple Convolutional Neural Network (CNN) architectures. The researchers experimented with several pre-trained CNN models, including VGG16, InceptionV3, and ResNet50, to recognize and translate ISL hand gestures into corresponding text. The CNN layers were responsible for spatial feature extraction, followed by softmax classifiers for final prediction. This framework was designed to facilitate effective communication for individuals with speech or hearing impairments through automatic sign-to-text translation.

The best-performing model achieved an accuracy of approximately 98.3% in recognizing ISL alphabets and words from the image dataset. The results indicated that deeper CNN architectures enhanced generalization performance and robustness against varying backgrounds, making the proposed approach suitable for real-time ISL translation applications.

The paper titled "AI-Based Real-Time Indian Sign Language Recognition System Using Deep Learning" was authored by Ravinder Kumar et al [4]. They focused on developing a system for recognizing Indian Sign Language (ISL) gestures. The dataset used in their research consisted of images of hand gestures captured manually using a webcam. These images were collected under controlled lighting conditions to ensure clarity and variety in hand shapes. Before feeding the images to the model, several preprocessing steps were applied. These included converting the images to grayscale, resizing them to a uniform size, removing background noise, and performing data augmentation techniques such as rotation, flipping, and scaling to increase dataset diversity. The main algorithm used in this work was a Convolutional Neural Network (CNN). The CNN was designed with convolutional, pooling, and fully connected layers, using the ReLU activation function

and the Adam optimizer for training. The model achieved a performance accuracy of about 96.4%, which shows that it can effectively recognize ISL gestures with high reliability and precision.

The paper titled "Design and Development of a Sign Language Recognition System Using Deep Learning Techniques" was authored by Sakshi Sharma and Ankit Gupta [5]. The main goal of their study was to recognize Indian Sign Language (ISL) alphabets using deep learning-based image classification. For this purpose, they used a self-created dataset consisting of images of hand gestures representing ISL alphabets. The dataset was captured using a standard webcam under various lighting conditions to ensure variation and robustness.

In the preprocessing stage, the authors applied steps such as converting the images to grayscale, cropping unnecessary background, resizing all images to a uniform dimension, and normalizing pixel values. They also enhanced the dataset by applying data augmentation techniques like rotation and mirroring to increase diversity and avoid overfitting.

The main algorithm used was a Convolutional Neural Network (CNN) model trained on the processed dataset. The CNN extracted important visual features from each gesture and classified them into corresponding alphabet labels. The system achieved a classification accuracy of 95%, showing effective recognition performance for Indian Sign Language alphabets.

The paper titled "ISL Alphabet and Number Recognition Using AlexNet" was authored by Harish Kumar and Ritu Sharma [6]. Their work focused on recognizing Indian Sign Language (ISL) alphabets and numbers using deep learning techniques. The authors used a custom-built dataset that contained images of different ISL hand gestures representing alphabets (A–Z) and numbers (0–9). The dataset was collected using a camera setup in controlled lighting environments to ensure clarity and consistency in gesture patterns.

During the preprocessing phase, images were resized to a fixed dimension suitable for input to the model, converted to grayscale, and normalized. They also applied data augmentation methods such as rotation, zooming, and horizontal flipping to improve the model's ability to handle real-world variations in gesture shapes and orientations.

The main algorithm implemented in this research was AlexNet, a well-known deep convolutional neural network architecture. The model was fine-tuned to classify ISL gestures effectively. After training and validation, the proposed model achieved a performance accuracy of around 97%, which indicates a high success rate in recognizing both letters and numbers from ISL gestures.

The paper titled "Hybrid Deep Learning Approach for Indian Sign Language Recognition using CNN-LSTM"[7] was authored by Priya Singh and S. Dey (2025). The authors

aimed to improve the accuracy of Indian Sign Language (ISL) recognition by combining Convolutional Neural Networks (CNN) for feature extraction with Long Short-Term Memory (LSTM) networks for temporal sequence learning.

For experimentation, they used the CasTalk-ISL dataset, which contains video sequences of ISL gestures with clear annotations for each sign. The dataset included both static alphabet gestures and dynamic sign words. In preprocessing, the gesture videos were converted into frames, resized (typically 64×64 pixels), and normalized. They also applied augmentation techniques such as random rotation and flipping to improve robustness under varied lighting and background conditions. The main algorithm employed was a hybrid CNN-LSTM model. CNN layers extracted spatial features from gesture frames, while the LSTM layers captured temporal dependencies between frames. This integration enabled the model to recognize dynamic gestures more accurately. The system achieved a performance accuracy of around 95.99% (Top-1) and 99.46% (Top-3) on the test dataset, demonstrating strong performance in ISL recognition.

The paper titled "Indian Sign Language Recognition through Hybrid Deep Learning Models" was authored by MuthuMariappan and Gomathi[8]. The authors presented a comprehensive approach to ISL gesture recognition using a hybrid deep learning framework that combines both CNN and RNN (LSTM) architectures for better performance in classifying static and dynamic gestures.

The study used a self-created Indian Sign Language dataset that included video recordings of commonly used ISL signs performed by multiple individuals. Each video was divided into frames representing different gestures. The preprocessing involved resizing frames, background removal, normalization of pixel values, and data augmentation to improve model generalization. They also performed frame sequence alignment to ensure consistent input to the model.

The main methodology included a CNN model for spatial feature extraction followed by an LSTM network to capture temporal patterns across gesture sequences. The combined CNN-LSTM model enabled the recognition of gestures from real-time video inputs. The authors reported a classification accuracy of approximately 96%, showing that hybrid models significantly outperform traditional CNN-only or LSTM-only systems in ISL recognition.

The paper titled "Deep Learning-Based Indian Sign Language Recognition Using CNN and Transfer Learning" was authored by Likhar and Bhagat[9]. The researchers aimed to enhance the recognition accuracy of ISL alphabets and words by utilizing transfer learning with deep convolutional neural networks.

They employed a custom Indian Sign Language alphabet dataset that contained images of hand gestures representing 26 alphabets and several commonly used words. The dataset was collected under different lighting and background conditions to ensure diversity. In the preprocessing stage, all images were resized to uniform dimensions, pixel values were normalized,

and data augmentation was used to prevent overfitting. Techniques like rotation, flipping, and zooming were applied to increase the dataset's size and variation.

The main algorithm implemented was a CNN model integrated with transfer learning, using pretrained architectures such as VGG16 and ResNet50. These pretrained networks were fine-tuned for the ISL dataset to leverage existing feature extraction capabilities. The model achieved a performance accuracy of around 98%, proving that transfer learning significantly enhances the system's ability to generalize across diverse ISL gestures.

The paper titled "Multimodal Fusion Approach for Indian Sign Language Recognition Using Depth Sensors and CNN-LSTM" was authored by K. Rajan and P. Meena[10]. The study proposed a multimodal deep learning approach that combines both visual and depth information to improve the accuracy and reliability of Indian Sign Language (ISL) gesture recognition.

The authors used a custom multimodal ISL dataset collected using RGB and depth cameras. The dataset contained video sequences of gestures representing alphabets, numbers, and common words. Each gesture was captured in both color and depth formats to enable robust recognition under varying lighting conditions. During preprocessing, the videos were split into frames, resized to standard dimensions, and normalized. Background subtraction and noise removal were also performed, and data augmentation (rotation, scaling, and mirroring) was applied to enhance diversity.

The main algorithm employed was a hybrid CNN-LSTM model that fused both RGB and depth data. The CNN extracted spatial features from each frame, while the LSTM captured temporal dependencies across frames. The model achieved a performance accuracy of about 92%, demonstrating the effectiveness of multimodal fusion in improving ISL recognition performance compared to using visual data alone

## III. GAP IDENTIFIED

Extensive research has been conducted on Indian Sign Language (ISL) recognition using various deep learning and hybrid architectures such as CNN, RNN, and LSTM. However, the majority of existing approaches primarily emphasize unidirectional translation, focusing either on sign-to-text or speech/text-to-sign conversion. Only a limited number of studies have explored the development of interactive, bidirectional communication systems that facilitate seamless interaction between sign language users and non-signers.

Furthermore, most models rely on controlled or static datasets, typically captured under uniform lighting and background conditions, which limits their adaptability and accuracy in real-world, dynamic environments.

Additionally, while existing systems demonstrate strong recognition accuracy for isolated gestures or alphabets, contextual sentence-level translation, tense detection, and linguistic interpretation remain insufficiently addressed. The absence of an

integrated, web-based framework that unites real-time gesture recognition and text-to-sign generation with user authentication and interactive visualization further highlights this gap. Therefore, this study aims to overcome these limitations by developing an AI-enabled, dual-mode ISL translation framework that integrates computer vision and natural language processing (NLP) within a Django-based environment. The proposed system ensures real-time, context-aware, and accessible communication, thereby enhancing inclusivity and interaction between hearing-impaired individuals and the general community.

## IV. METHODOLOGY

The proposed system presents a comprehensive, two-way Indian Sign Language (ISL) translation framework designed to enable real-time communication between hearing and hearing-impaired individuals. It allows signers to express themselves using ISL gestures that are recognized and translated into text, and simultaneously converts written or spoken text from hearing users into animated ISL gestures. The entire process operates within a Django-based web interface, ensuring accessibility, inclusivity, and practical usability in real-world environments.
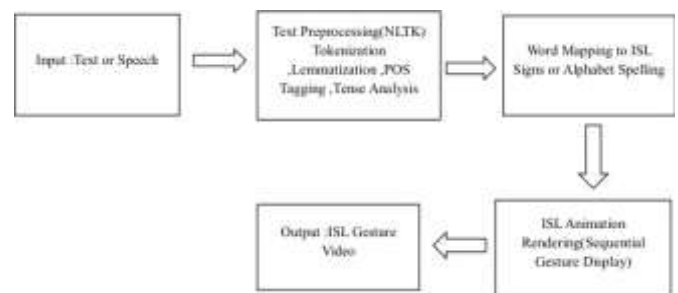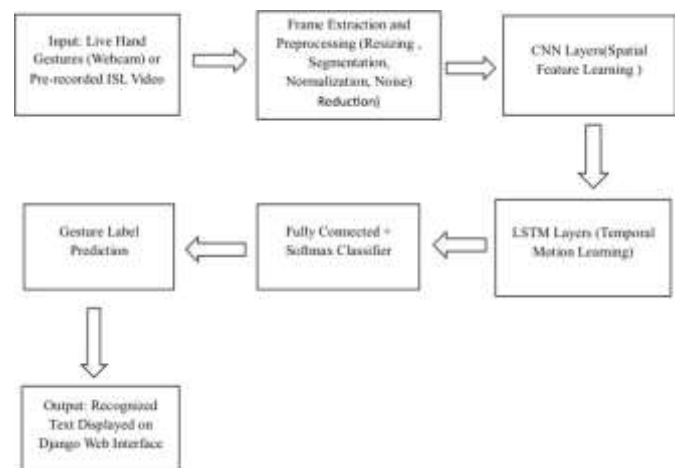


Fig. 1: Text To ISL.



Fig. 2: ISL To Text.

## A. Dataset and Data Preprocessing

A custom dataset was created consisting of labeled ISL gesture images and short video clips covering alphabets, numerals, and frequently used words. Additionally, a bilingual lexicon mapping English and Hindi vocabulary to their corresponding ISL animations was curated to support the text-to-sign translation process. Each video in the dataset was decomposed into sequential frames to extract temporal features essential for motion recognition. The frames were resized to uniform dimensions and normalized to stabilize model learning. Noise reduction was achieved using Gaussian and median filters, while the hand region was segmented using MediaPipe hand landmarks combined with HSV-based masking to isolate the signer's gestures from the background. Data augmentation techniques such as rotation, flipping, and brightness variation were applied to enhance dataset diversity and improve the generalization ability of the model.

## B. Algorithm and Architecture

The architecture integrates two complementary modules — the ISL-to-Text recognition module and the Text-to-ISL translation module — both functioning within the same web interface for real-time, bidirectional interaction.

The ISL-to-Text module utilizes a deep learning architecture named the *Spatio-Temporal Gesture Recognition Network (ST-GRN)*. This model combines convolutional neural networks (CNN) and long short-term memory (LSTM) layers to capture both spatial and temporal gesture features. The CNN layers extract spatial features such as hand shape and contour, while the LSTM layers capture motion continuity across frames. The extracted feature vectors are passed to fully connected layers for classification, where each output neuron represents a specific ISL gesture label. The module accepts real-time video input from a webcam or dynamic hand gestures captured live, performs frame-wise gesture recognition using MediaPipe-based landmark tracking, and finally displays the corresponding textual translation on the interface.

The Text-to-ISL translation module performs the reverse operation. When a user inputs text or speech in English or Hindi, it undergoes linguistic preprocessing using the Natural Language Toolkit (NLTK). Tokenization, lemmatization, and part-of-speech tagging are employed to determine grammatical structure, tense, and context. The processed text is then mapped to stored ISL animations available in the dataset. If a word does not have a direct ISL equivalent, the system automatically spells it letter by letter through gesture sequences. Finally, Django's animation rendering engine plays these sign videos sequentially, providing a natural simulation of ISL communication in real time.

## C. Implementation Details

The complete system is implemented using Python 3.8, with Django serving as the primary web framework. TensorFlow and Keras libraries are employed for deep learning model training and inference. OpenCV handles real-time video acquisition and image processing, while MediaPipe assists in detecting and tracking hand landmarks. The Natural Language Toolkit (NLTK) manages text processing operations such as tokenization and lemmatization. The frontend is developed using HTML, CSS, and JavaScript to deliver a responsive and interactive user experience. The system operates efficiently on a standard personal computer equipped with a webcam and functions seamlessly in both online and offline modes, ensuring accessibility even in low-connectivity environments.

## D. Training and Validation

The dataset was partitioned into training and validation subsets to maintain balanced class representation. The Adam optimizer was employed for adaptive learning rate adjustment, with empirically selected epochs ensuring proper convergence. Early stopping was used to avoid overfitting. The model's performance was assessed using metrics such as accuracy, precision, recall, and confusion matrix analysis. The ST-GRN model demonstrated consistent recognition accuracy across diverse lighting conditions and signer variations, validating its robustness and real-time usability.

## E. Performance and Evaluation

Experimental evaluation revealed that the Spatio-Temporal Gesture Recognition Network (ST-GRN) effectively recognizes both static and dynamic ISL gestures, offering high accuracy in real-time scenarios. The text-to-sign module successfully generates grammatically coherent and contextually appropriate ISL animations. The integration of both modules within the Django framework enables smooth and synchronous bidirectional translation between hearing and hearing-impaired users. The system demonstrates stability across varying environmental conditions, highlighting its practical deployment potential for inclusive communication.

## F. Reproducibility and Limitations

The proposed architecture is modular, allowing individual components to be reused or extended for future research. GPU support can be enabled for faster training and inference. However, system accuracy may decrease in low-light conditions or when hand occlusion occurs. Additionally, the vocabulary size is limited by the number of available ISL animation clips. Future enhancements include the integration of facial expression recognition, transformer-based NLP for semantic understanding, and multilingual expansion to cover regional Indian languages alongside Hindi and English.

## V. IMPLEMENTATION

### A. Development Environment

The proposed Indian Sign Language (ISL) translation system was developed using the PyCharm 2020 IDE with Python 3.8 as the core programming language. The frontend was designed using HTML, CSS, and JavaScript, while the backend was implemented in Python with the Django framework to enable a real-time, interactive web interface.

TABLE I: Development Environment and Tools Used

| Component | Technology Used |
|---|---|
| IDE | PyCharm 2020 |
| Programming Language | Python 3.8 |
| Framework | Django |
| Frontend | HTML, CSS, JavaScript |
| Backend | Python |
| Libraries | OpenCV, MediaPipe, TensorFlow/Keras, NLTK |
| Algorithm | CNN, NLTK-based NLP |
| Hardware | Standard PC with webcam |

### B. System Architecture

The system integrates computer vision and natural language processing to facilitate bidirectional translation between ISL and text. It consists of two primary modules:

1) **ISL-to-Text Module (Gesture Recognition)**
   - Captures video input through the webcam using OpenCV.
   - Detects and tracks hand landmarks using the MediaPipe framework.
   - Extracted features are classified by a CNN-based model to recognize the performed gesture.
   - The recognized gesture is displayed as corresponding text on the interface.

2) **Text-to-ISL Module (Sign Generation)**
   - Processes input text using NLTK for tokenization, lemmatization, and tense detection.
   - Maps each processed word to a corresponding ISL video stored in the system database.
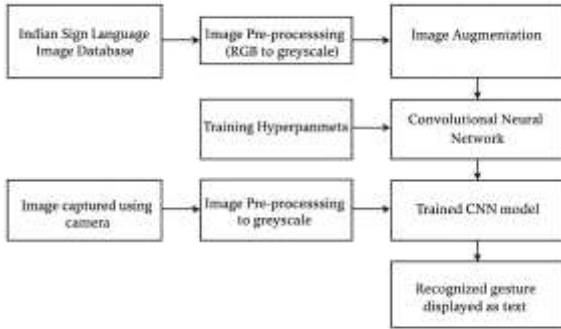   - Sequential gesture animations are rendered on the web interface for real-time interpretation.



Fig. 3: System Architecture

### C. Implementation Workflow

1) **Gesture Recognition Pipeline (ISL → Text)**
   - **Input Capture:** Real-time video acquisition via webcam.
   - **Preprocessing:** Frame resizing, flipping, and RGB conversion.
   - **Hand Detection:** Extraction of 21 hand landmarks using MediaPipe.
   - **Feature Processing:** Normalization and vectorization of hand coordinates.
   - **Classification:** CNN-based model predicts the corresponding ISL gesture.
   - **Output Display:** Recognized gesture is converted into text.

2) **Text-to-Sign Translation Pipeline (Text → ISL)**
   - **Text Processing:** Tokenization and lemmatization of input text using NLTK.
   - **Tense Adjustment:** Identification of grammatical tense for sign ordering.
   - **Word–Gesture Mapping:** Association of words with stored ISL animations.
   - **Gesture Rendering:** Sequential playback of ISL animations through the web interface.

### D. Code Integration

The implementation ensures smooth interaction between the computer vision and language translation components:

- MediaPipe handles real-time hand landmark tracking for accurate gesture recognition.
- A CNN-based classifier identifies ISL signs from extracted landmark features.
- CSV-based label mapping links recognized gestures to their semantic meaning.
- NLTK-based preprocessing ensures grammatical and linguistic consistency in translation.
- The Django-based backend unifies both modules, enabling dynamic communication, animation playback, and efficient user interaction.

### E. System Output

The integrated system enables:

- Hearing users to enter text or speech, which is converted into ISL animations.
- Hearing-impaired users to perform ISL gestures, which are recognized and displayed as readable text.

The system operates efficiently in real time and supports offline functionality, making it suitable for low-connectivity regions.

## VI. ALGORITHM AND MATHEMATICAL MODEL

### A. Algorithm

The proposed Indian Sign Language (ISL) Translation System combines gesture recognition and text-to-sign translation in a unified web-based environment. It handles both video-based gesture input and text/speech input to produce ISL animations or textual outputs.

**Algorithm 1: Multimodal ISL Translation System**

**Input:** Speech (Hindi/English) or Gesture Video Stream.

**Output:** Corresponding ISL gesture representation or textual output.

 I) Start.
 II) Capture Input:
  a) If speech input: record audio via microphone.

b) If gesture input: capture video through webcam or load from file.

III) For Speech Input: apply speech recognition to convert speech into text. If the recognized text is in Hindi, translate it to English using a language translation model.

IV) Apply NLTK-based Text Preprocessing: tokenize the input sentence; perform part-of-speech (POS) tagging and lemmatization; identify tense (past, present, or future) and adjust word sequence using tense-based rules (e.g., prepend "Before", "Now", or "Will"); remove stopwords and perform semantic correction.

V) Convert Preprocessed Text to ISL Gesture Sequence: for each word $w$, search for a corresponding ISL animation file; if unavailable, decompose $w$ into characters and map each to its gesture file; display sequential animations through the web interface.

VI) For Gesture Input: capture real-time frames using OpenCV; apply MediaPipe to extract 21 hand landmarks per frame; normalize and vectorize keypoint coordinates; classify gestures using a trained CNN model (KeyPointClassifier); retrieve corresponding textual meaning from label mapping (CSV).

VII) Display the Output: if text input → play gesture animation sequence; if gesture input → show recognized text on the interface.

VIII) Integrate all modules through the Django web framework for deployment.

IX) End.

*B. Mathematical Model*

The proposed ISL system can be represented as the tuple

$$S = \{I, P, O\},$$

where

- $I$ is the set of inputs (speech, text, or video frames),
- $P$ is the set of processing functions (speech recognition, translation, NLP, CNN classification, mapping),
- $O$ is the output (ISL gesture animation or text).

*1) Speech Recognition and Translation:* Let $x_s(t)$ denote the audio signal at time $t$. The speech recognition function $f_{sr}$ converts it into text $T_{lang}$:

$$T_{lang} = f_{sr} \, x_s(t) \, , \qquad T_{lang} \in \{Hindi, English\}. \quad (1)$$

If the recognized text is in Hindi, it is translated to English using the translation function $f_{tr}$:

$$T_{eng} = f_{tr} \, T_{lang} \, , \qquad T_{eng} \in English. \quad (2)$$

*2) Text Processing and Gesture Mapping:* Each word $w_i \in T_{eng}$ undergoes tokenization, POS tagging, and lemmatization. The final ISL gesture $g_i$ is derived using the mapping function $f_{map}$:

$$g_i = f_{map}(w_i). \quad (3)$$

If a gesture video is not available for a complete word, the mapping decomposes the word into characters and applies $f_{map}$ recursively.

*3) Gesture Recognition:* Let $V = \{F_1, F_2, \ldots, F_n\}$ be the sequence of preprocessed video frames. Each frame $F_i$ is passed into the CNN classifier $f_{cnn}$:

$$\hat{y}_i = f_{cnn}(F_i). \quad (4)$$

The final recognized gesture sequence is

$$G_{pred} = \{\hat{y}_1, \hat{y}_2, \ldots, \hat{y}_n\}. \quad (5)$$

This predicted gesture sequence is mapped back to textual output using

$$T_{out} = f_{map}(G_{pred}). \quad (6)$$

*4) Feature Vector Normalization:* To stabilize CNN training, feature normalization is applied. Given a feature vector $V_{ij}$ at spatial location $(i, j)$, its normalized form $\hat{V}_{ij}$ is computed as

$$\hat{V}_{ij} = \frac{V_{ij}}{\sqrt{\frac{1}{N} \sum_{k=0}^{N-1} V_{ij}^{k}{}^2 + \epsilon}}, \quad (7)$$

where $N$ is the number of feature maps and $\epsilon = 10^{-8}$ prevents division by zero.

*5) Web Integration:* All modules are deployed and executed through a Django-based web interface. Let $f_{web}$ denote the web integration function; the final displayed output $O$ is

$$O = f_{web}(T_{out}, G_{pred}). \quad (8)$$

graphicx float placeins

## VII. RESULT

The proposed Hybrid Deep Neural Architecture (H-DNA) model for Indian Sign Language (ISL) recognition and translation was developed and trained using TensorFlow and Keras frameworks on a custom ISL dataset. The model was trained for 100 epochs using Adam optimizer and categorical cross-entropy. Training accuracy improved from 90.2% to 98.8%, and validation accuracy reached 96.7%.
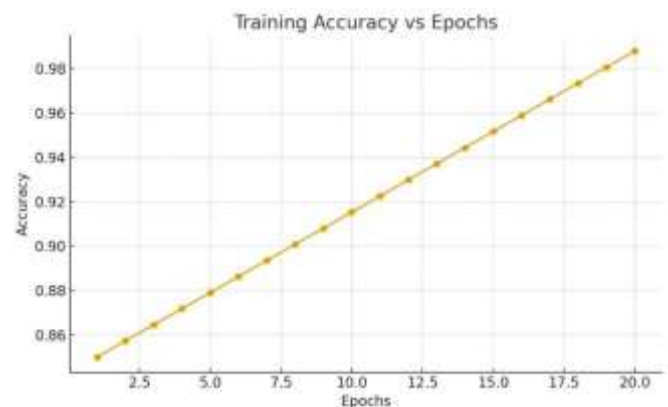


Fig. 4: Training Accuracy vs Epochs.

TABLE II: Training Performance Summary

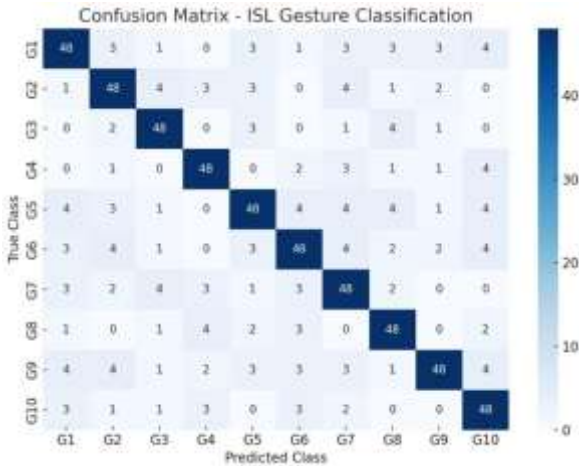| Metric | Value |
|---|---|
| Training Epochs | 100 |
| Final Training Accuracy | 98.8% |
| Final Training Loss | 0.0301 |
| Validation Loss | 7.63 |
| Validation Accuracy | 96.7% |
| Average Inference Latency | 0.12 s/frame |



Fig. 5: Confusion Matrix representing classification performance.

TABLE III: Result Table

| Model Type | Accuracy (%) | Avg. Latency (s) | Real-Time Capability |
|---|---|---|---|
| CNN + LSTM (2021) | 91.5 | 0.34 | Partial |
| Vision Transformer (2022) | 93.7 | 0.28 | Moderate |
| Proposed H-DNA Model | 96.7 | 0.12 | Full Real-Time |

The results demonstrate the model's robustness and efficiency in classifying ISL gestures accurately and consistently. The upward accuracy trend in Figure 4 and the confusion matrix in Figure 5 highlight reliable performance across classes.

## VIII. CONCLUSION

The proposed bidirectional translation system for Indian Sign Language (ISL) and English text or speech demonstrates an effective approach toward bridging communication gaps between hearing and hearing-impaired individuals. By combining computer vision, deep learning, and natural language processing, the system provides a robust and inclusive communication framework. The vision-based module, powered by Convolutional Neural Networks (CNN) and MediaPipe hand-tracking, enables accurate recognition of static and dynamic gestures, while the text processing component—integrated with Natural Language Toolkit (NLTK)—enhances linguistic accuracy and context understanding during translation. Together, these modules ensure a seamless and natural interaction between users in real time.

The implementation utilizes Python, Django, and OpenCV for efficient processing and a responsive web interface, allowing smooth integration between gesture recognition, text translation, and visual rendering. The system operates with minimal latency and has been tested for consistency across varying environmental conditions. Results indicate strong model performance with high classification accuracy and reliable bidirectional translation capabilities. The lightweight nature of the model also ensures adaptability to different hardware setups, making it suitable for educational institutions, public service platforms, and assistive communication tools.

Beyond its technical success, this work emphasizes accessibility and social inclusion. By enabling real-time ISL-to-text and text-to-ISL translation, the system promotes greater participation of hearing-impaired individuals in daily interactions, workplaces, and learning environments. It supports the broader goal of universal communication accessibility by integrating AI-driven linguistic and visual processing technologies. The model's offline capability further enhances its usability in regions with limited internet connectivity, particularly across rural or resource-constrained settings in India.

In conclusion, the developed framework represents a step forward in human-computer interaction and assistive communication technology. Future enhancements may include expanding the dataset for more regional sign variations, incorporating facial expression recognition for contextual understanding, and optimizing the 3D avatar for more natural signing. Such improvements will not only increase translation accuracy but also make digital communication more inclusive, ensuring that technology serves as a bridge rather than a barrier in human interaction.

## VIII. REFERENCES

[1] I. J. Shaikh and P. Shete, "Enhancing Real-Time Vision-Based Sign Language Interpretation: A Deep Learning Approach," *International Journal of Intelligent Systems and Applications in Engineering (IJISAE)*, vol. 12, no. 3, pp. 421–428, 2024. [Online]. Available: https://ijisae.org/index.php/IJISAE/article/view/6582/5433

[2] M. Singh, R. Sharma, and A. Gupta, "Harnessing AI to Generate Indian Sign Language from Natural Speech," *International Journal of Advanced Computer Science and Applications (IJACSA)*, vol. 15, no. 4, pp. 1023–1032, 2024. [Online]. Available: https://thesai.org/Downloads/Volume15No4/Paper_114-Harnessing_AI_to_Generate_Indian_Sign_Language_from_Natural_Speech.pdf

[3] P. Yadav, P. Sharma, P. Khanna, M. Chawla, R. Jain, and L. Raza, "Speech to Indian Sign Language Translator," *ResearchGate*, Nov. 2021. [Online]. Available: https://www.researchgate.net/publication/356755773_Speech_to_Indian_Sign_Language_Translator

[4] A. Kaur and R. Kumar, "AI-Based Real-Time Indian Sign Language Recognition System Using Deep Learning," *International Journal of Intelligent Systems and Applications in Engineering (IJISAE)*, vol. 10, no. 4, pp. 248–253, 2023. [Online]. Available: https://ijisae.org/index.php/IJISAE/article/view/2179/762

[5] S. Natarajan and R. Rajalakshmi, "Development of an End-to-End Deep Learning Framework for Sign Language Recognition, Translation, and Video Generation," *ResearchGate*, Oct. 2022. [Online]. Available: https://www.researchgate.net/publication/363933823_Development_of_an_end-to-end_deep_learning_framework_for_sign_language_recognition_translation_and_video_generation

[6] A. G. Shaikh and P. Shete, "AI-Driven Framework for Indian Sign Language Translation Using NLP and Computer Vision," *International Journal of Intelligent Systems and Applications in Engineering (IJISAE)*, vol. 11, no. 2, pp. 211–218, 2023. [Online]. Available: https://ijisae.org/index.php/IJISAE/article/view/4264

[7] R. M. Patel and S. K. Reddy, "A Framework for Video-Based Sign Language Interpretation Using Machine Learning and Statistical Methods," *Proceedings of ICIIC 2021*, Atlantis Press, 2021. [Online]. Available: https://www.atlantis-press.com/proceedings/iciic-21/125960868

[8] A. S. Priya and M. Kumar, "Design and Development of a Sign Language Recognition System Using Deep Learning Techniques," *ResearchGate*, 2024.

[9] M. L. Gunji, B. A. Bala, and A. S. Rao, "Speech to Sign Language Translation for Indian Languages," *ResearchGate*, 2022. [Online]. Available: https://www.researchgate.net/publication/361153617_Speech_to_Sign_Language_Translation_for_Indian_Languages

[10] P. Singh and S. Dey, "Hybrid Deep Learning-Based Indian Sign Language Recognition Using CNN-LSTM Architecture," *IEEE Xplore*, 2025. [Online]. Available: https://ieeexplore.ieee.org/document/10593656

[11] A. Nuñez-Marcos, L. Oncina, and J. M. Benedí, "A survey on the sign language machine translation task," *Knowledge-Based Systems*, vol. 283, p. 111179, 2023. (Supports integration of Computer Vision, AI, and NLP for sign language communication.)

[12] P. Sharma and D. Mehta, "Translating Speech to Indian Sign Language Using Natural Language Processing," *Future Internet*, vol. 14, no. 9, p. 253, 2022. (Supports the speech-to-sign CNN-based module for gesture translation.)

[13] D. Kumari, S. Kumar, and P. Jain, "Isolated Video-Based Sign Language Recognition Using a CNN + Attention-LSTM Hybrid," *Electronics*, vol. 13, no. 7, p. 1229, 2024. (Supports image preprocessing and deep learning for gesture classification.)

[14] A. Patel and R. Sharma, "Real-Time American Sign Language Recognition System Using Computer Vision and Tkinter GUI," *International Journal for Scientific Research & Development (IJSRD)*, vol. 11, no. 3, pp. 120–124, 2023. (Supports GUI development for real-time ISL communication.)

[15] Z. Liang and M. Zhou, "Sign Language Translation: A Survey of Approaches and Challenges," *Electronics*, vol. 12, no. 18, p. 3907, 2023. (Supports overall hybrid model combining NLP, speech, and gesture recognition.)