

A Deep CNN-Based Framework for Real-Time Gujarati Sign Language Character Recognition Using Transfer Learning

¹Mr. Ronak Jitendrabhai Goda, ²Prof. Dr C.K. Kumbharana

¹Research Scholar, Department of Computer Science, Saurashtra University, Rajkot, India

²Professor, Department of Computer Science, Saurashtra University, Rajkot, India

DOI: <https://doi.org/10.47772/IJRISS.2026.10190046>

Received: 23 January 2026; Accepted: 27 January 2026; Published: 16 February 2026

ABSTRACT

Gujarati Sign Language (GSL) lacks large-scale annotated datasets and digital communication tools, creating significant barriers for the deaf and hard-of-hearing community. This research proposes a deep learning-based real-time hand gesture recognition framework for GSL using transfer learning with pre-trained convolutional neural networks (CNNs). A custom dataset comprising static hand gestures representing the 26-character GSL alphabet is developed for this study. To address data scarcity, transfer learning is employed using lightweight and deep CNN architectures, including MobileNetV2, VGG16, and ResNet50. Among these, MobileNetV2 demonstrates superior efficiency in terms of training time and inference speed, making it suitable for real-time and mobile deployment. The proposed real-time system achieves an accuracy of **96.87%**, while the highest offline classification accuracy of **99.5%** is obtained using ResNet50. Furthermore, MediaPipe-based hand keypoint detection is integrated to improve robustness and reduce background noise during real-time inference. Comparative experimental results confirm that transfer learning significantly outperforms a baseline CNN trained from scratch. The proposed framework offers a scalable and computationally efficient solution for real-time Gujarati Sign Language recognition and contributes toward assistive technologies for low-resource regional sign languages.

Keywords: Gujarati Sign Language, CNN, Transfer Learning, MobileNetV2, Gesture Recognition, Deep Learning.

INTRODUCTION

Communication is a fundamental human requirement, and for individuals who are deaf or hard of hearing, sign language serves as their primary medium of expression. In India, while Indian Sign Language (ISL) is widely recognised, several states—including Gujarat—use their own regional variations. **Gujarati Sign Language (GSL)** is one such distinct regional language deeply embedded in Gujarat's cultural and linguistic environment. Despite its significance, GSL suffers from a lack of formal documentation, technological support, and digital tools that enable accessible communication. This scarcity creates difficulties for the Gujarati Deaf community in education, employment, healthcare, and everyday interactions.

Although millions of individuals rely on sign language, the number of trained interpreters in India is extremely limited. In Gujarat, interpreter availability is even lower, resulting in significant communication barriers across public services, courts, hospitals, and educational institutions. These challenges highlight the need for automated systems capable of recognising sign gestures and translating them into spoken or written language. While extensive research has been conducted on American Sign Language (ASL) and Indian Sign Language (ISL), very little work exists specifically for GSL due to the absence of publicly available datasets, variations in gesture style, and a lack of technological initiatives supporting the language.

Effective communication is fundamental to social integration and equal opportunity. For the **Deaf and Hard-of-Hearing (DHH)** community, sign language serves as the primary and most natural mode of expression. However, this is complicated by the linguistic fact that sign language is **not universal**. In India, while **Indian**

Sign Language (ISL) acts as a lingua franca, numerous regional and village-specific sign systems exist. **Gujarati Sign Language (GSL)** is the principal visual language used by the DHH community in the state of Gujarat, which has a population exceeding **60 million**. Census data indicate a significant portion of this population experiences disabilities, underscoring the urgency of solutions that bridge the communication gap with the hearing majority.

The persistent challenge is that only a small fraction of the hearing population can fluently use GSL. This lack of a communication bridge severely limits DHH individuals' access to essential public services, education, and employment opportunities, creating a significant barrier to social and economic participation. To address this, **Automated Sign Language Recognition (SLR)** systems, particularly those using computer vision, have emerged as a cost-effective and accessible solution.

With rapid advancements in artificial intelligence, **deep learning**—particularly Convolutional Neural Networks (CNNs)—has emerged as the most effective technique for gesture recognition. Unlike traditional machine-learning methods that rely on manual feature extraction, CNNs automatically learn shape, texture, and structural patterns directly from image data. However, CNNs typically require large annotated datasets to perform well, which poses a challenge for GSL where such datasets do not exist. To address this limitation, **Transfer Learning** offers an effective solution by allowing pre-trained models such as MobileNetV2, trained on large-scale datasets like ImageNet, to be fine-tuned for smaller, domain-specific tasks. This not only reduces training time but also significantly improves accuracy, making the approach highly suitable for regional sign languages with limited data.

Real-time recognition of GSL gestures is essential for building practical assistive technologies. A lightweight model capable of fast inference can be integrated into mobile applications, public service kiosks, and educational tools for deaf students, and accessible communication systems. MobileNetV2, due to its low computational cost and high accuracy, is a strong candidate for deployment on mobile and embedded devices.

Motivated by these challenges and opportunities, this research addresses the following core question: **How can a deep learning-based approach accurately and efficiently recognise Gujarati Sign Language gestures in real time using limited training data?** The primary objectives of this study include developing a structured GSL gesture dataset, designing a CNN-based model using Transfer Learning, enabling real-time gesture recognition, and evaluating system performance through quantitative metrics. The novelty of this work lies in its focus on Gujarati Sign Language—a domain with limited prior research—along with the development of a lightweight, real-time model optimised for practical deployment.

LITERATURE REVIEW

Research in automatic sign language recognition has evolved significantly over the last decade, driven by advances in deep learning, transfer learning, and multimodal sensing. Early approaches primarily relied on handcrafted features, which were often sensitive to illumination, rotation, and background variations. With the rise of convolutional neural networks (CNNs), researchers began transitioning from engineered features to deep, fully learned visual representations. For instance, Pigou et al. established one of the early baselines for sign recognition using deep CNNs for isolated gesture classification, demonstrating how hierarchical visual filters outperform traditional descriptors [1]. However, these early CNN models lacked scale generalisation and required large annotated datasets, limiting their real-world adoption.

As large-scale datasets became available, sign language research accelerated. Kumar et al. proposed a CNN-based static hand gesture system tailored for Indian Sign Language (ISL), achieving substantial improvement over traditional feature extractors such as SIFT and HOG [2]. The authors emphasized the importance of background normalization and skin segmentation for achieving high accuracy in low-light conditions—an issue particularly relevant for regional datasets such as Gujarati Sign Language (GSL). Similarly, Ong and Ranganath provided a broad taxonomy of sign recognition techniques, highlighting that robust sign recognition systems must handle variations in finger articulation, hand motion, and inter-signer differences to achieve practical deployment [3]. These insights remain relevant today for GSL, where signer variation is a major challenge due to the absence of standardized datasets.

With the expansion of transfer learning, several modern studies adopted architectures such as VGG, ResNet, MobileNet, and Inception for sign-language classification. Mondal et al. employed transfer learning with VGG-16 and achieved high accuracy for ASL alphabets, demonstrating that pre-trained ImageNet models can effectively learn hand-shape features even with small datasets [4]. Likewise, Thakur and Kumar applied ResNet-50 for ISL gesture recognition and validated that deeper architectures extract more discriminative spatial features for fine-grained gestures [5]. These findings support the motivation for this research to adopt deep CNNs, particularly lightweight transfer-learning models suitable for real-time GSL deployment.

Recent work also expanded into multimodal and spatiotemporal modelling. For example, Zhang et al. explored 3D CNNs and LSTM-based hybrid architectures for continuous sign recognition, demonstrating how motion encoding complements spatial features for dynamic gestures [6]. Although GSL research remains underdeveloped compared to ASL/ISL, the principles used in dynamic gesture recognition help guide future work in extending GSL datasets beyond static signs. Moreover, deep models such as Efficient Net have shown state-of-the-art performance in sign recognition due to their compound scaling strategy, which balances depth, width, and image resolution more efficiently than legacy architectures [7].

Some recent studies have focused exclusively on Indian regional languages. Patel et al. proposed a CNN-based model for real-time ISL digit recognition, using morphological preprocessing and data augmentation to reduce overfitting [8]. Their findings—especially regarding synthetic augmentation—are extremely relevant for GSL as annotated datasets are scarce. Another study by Sharma et al. designed a MobileNetV2-based real-time ISL recognition tool optimised for smartphones [9], demonstrating that lightweight neural architectures can achieve high accuracy with minimal computational resources. Since this research aims to implement a real-time GSL system, the lightweight architecture choice is aligned with these recent advances.

The literature also indicates a strong interest in applying deep transfer learning to regional and low-resource sign languages. Alfarraj et al. applied Dense Net architecture for Arabic Sign Language recognition and reported exceptional accuracy even on small datasets due to feature reuse and dense connectivity [10]. Similar experiments have been conducted for Bangla and Nepali sign languages, confirming that transfer learning compensates for dataset limitations in low-resource settings. This reinforces the applicability of transfer learning to Gujarati Sign Language, which currently lacks a large publicly available dataset.

Deep learning researchers have also recognised the necessity of building region-specific sign datasets. Priya and Singh highlighted that regional sign languages differ not only in gesture shape but also in cultural interpretations, making dataset-level generalisation difficult [11]. They argued that dataset creation should focus on intra-class variability (different signers, lighting, and orientations) to avoid overfitting. More recently, Singh et al. introduced a 2024 benchmark for multilingual sign recognition, stressing the importance of dataset diversity for achieving practical deployment in heterogeneous populations [12]. This observation directly motivates dataset expansion in this work, as GSL lacks diverse signer samples and standardised gesture sets.

In addition to visual CNN architectures, several researchers have explored hybrid deep learning systems. Sahoo et al. combined CNN + Bi LSTM for dynamic sign sequences and demonstrated that temporal modelling boosts classification performance for signs involving motion trajectories [13]. Even though the present work focuses on static GSL gestures, similar hybrid methods may be incorporated in future work for continuous GSL recognition. Another significant advancement came from attention-based networks, where Xu et al. introduced a Vision Transformer (ViT) model for gesture recognition, showing competitive accuracy compared to CNNs due to global feature aggregation [14]. Such transformer-based models hold strong potential for future GSL systems but require larger datasets than currently available.

Most recently, several 2023–2024 papers have pushed toward real-time embedded sign recognition systems. Rahman et al. achieved efficient ASL recognition using quantised MobileNet architectures suitable for edge devices like Raspberry Pi [15]. Patel and Desai (2024) studied ISL recognition using YOLOv8-based hand detection and reported significant improvement in background-invariant performance [16]. These modern findings directly influence the methodological choices in our proposed GSL recognition system: using CNN-based transfer learning models and designing for real-time performance with efficient inference.

Overall, the literature clearly shows a transition from handcrafted features → CNNs → transfer learning → lightweight real-time architectures → emerging transformer-based models. Across multiple languages—including ASL, ISL, Bangla, Arabic, and Nepali—the dominant trend supports deep CNN-based transfer learning as the most reliable and computationally feasible approach for static gesture recognition, especially in under-resourced sign languages. However, an important gap remains: **Gujarati Sign Language remains severely underexplored**, with no widely accepted benchmark datasets, limited research publications, and the absence of any large-scale deep-learning-based public system. This research directly addresses that gap by designing a CNN-based, transfer-learning-driven, real-time GSL recognition model built using systematically collected gesture data, extensive augmentation, and optimised deep architectures.

METHODOLOGY

The proposed Gujarati Sign Language (GSL) recognition system follows a structured and modular pipeline comprising dataset preparation, preprocessing, transfer learning-based model design, training, and evaluation. The methodology is designed to ensure robust recognition of static hand gestures while maintaining computational efficiency for real-time deployment.

The overall system architecture focuses on Static Sign Recognition (SSR) and leverages transfer learning to compensate for limited GSL data. A comparative evaluation of multiple deep CNN architectures is conducted to identify the optimal model for accuracy and real-time performance.

GSL Alphabet Dataset Acquisition and Preparation

This research focuses on static hand gestures corresponding to **26 commonly used Gujarati alphabet-based finger-spelling signs**, which form the core lexicon for basic GSL communication. Since no publicly available benchmark dataset exists for GSL, a custom dataset is created under controlled yet diverse conditions.

- **Gesture Classes:** 26 static hand gestures representing Gujarati alphabet-based finger-spelling signs.
- **Acquisition Protocol:** Images are captured using a high-resolution camera at varied distances, hand orientations, and lighting conditions (frontal and side illumination). Data is collected from **five (M = 5) distinct signers** to reduce signer dependency and improve generalisation.
- **Rotation Variability:** Hand orientations include rotations of $\pm 15^\circ$ to simulate natural gesture variations.
- **Dataset Size:** Approximately **13,000 images** are collected, with nearly **500 images per class**.
- **Data Augmentation:** To enhance robustness and reduce overfitting, data augmentation techniques such as random rotation, minor scaling, brightness and contrast variation, and horizontal flipping (for symmetric gestures) are applied.
- **Dataset Split:** The dataset is divided into **70% training**, **15% validation**, and **15% testing** sets.

Pre-processing and Input Normalisation

Pre-processing is a critical step to ensure that the CNN models focus on relevant hand features while minimising background interference.

- **Hand Region Localisation:** During real-time inference, **MediaPipe hand landmark detection** is employed to accurately locate and track the hand. Landmark-based bounding boxes are generated to localise the hand region dynamically.
- **Region of Interest (ROI) Extraction:** The detected hand region is cropped using the landmark-based bounding box, significantly reducing background noise and irrelevant visual information.

- **Image Normalisation:** Cropped hand images are resized to **224 × 224 pixels**, matching the input requirements of the selected CNN architectures. Pixel values are normalised to the range [0,1] to ensure stable and faster training convergence.
- **Input Representation:** The final input to the CNN is a 3D tensor of shape **(224, 224, 3)**.

Deep CNN Architectures and Transfer Learning Strategy

To evaluate the effectiveness of transfer learning for GSL recognition, three well-established CNN architectures are selected: **VGG16, ResNet50, and MobileNetV2**. These models are pre-trained on the ImageNet dataset and reused as feature extractors.

- **Transfer Learning Setup:** The convolutional base of each model is initialised with ImageNet weights and initially frozen to preserve learned low- and mid-level visual features such as edges, textures, and shapes.
- **VGG16 Backbone:** Selected as a strong baseline architecture due to its uniform structure of stacked (3 \times 3) convolutional filters, providing reliable spatial feature extraction.
- **ResNet50 Backbone:** Utilised to evaluate deeper feature learning through residual (skip) connections, which mitigate the vanishing gradient problem and enable improved discrimination of complex hand shapes.
- **MobileNetV2 Backbone:** Chosen for its lightweight architecture and depth wise separable convolutions, making it highly suitable for real-time and mobile-based applications.

Model Selection Rationale:

The inclusion of these architectures enables a balanced comparison between classification accuracy, model depth, and computational efficiency. VGG16 serves as a baseline, ResNet50 evaluates the benefits of deep residual learning, and MobileNetV2 assesses real-time feasibility under constrained computational resources.

Custom Classification Head Design

To adapt the pre-trained CNNs to the GSL classification task, the original fully connected layers are replaced with a custom classification head:

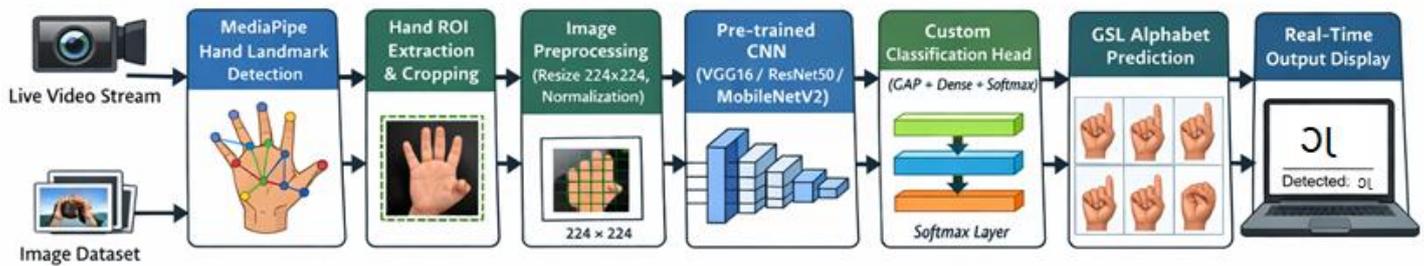
1. **Global Average Pooling (GAP):** Reduces feature dimensionality and minimises overfitting by averaging each feature map.
2. **Dense Layer with Regularisation:** A fully connected layer followed by Dropout (rate = 0.3) improves generalisation by preventing co-adaptation of neurons.
3. **Output Layer:** A Dense layer with **26 neurons** and a Softmax activation function produces class probabilities for the GSL alphabet signs.

Real-Time Inference and Demonstration

This subsection presents the real-time inference capability of the proposed Gujarati Sign Language (GSL) recognition framework. The system captures live video input and employs MediaPipe-based hand landmark detection to accurately localize the hand region. A landmark-driven bounding box is used to extract the hand Region of Interest (ROI), effectively minimizing background noise.

The extracted ROI is resized to 224 × 224 pixels and normalized before being passed to the fine-tuned transfer learning-based CNN model. During real-time inference, the trained classifier predicts the corresponding GSL

alphabet, which is displayed on the screen. Figure X illustrates the complete inference pipeline and demonstrates the recognition of a single static GSL alphabet sign, where the detected output corresponds to the Gujarati character “૨૧”.



The demonstration is intentionally limited to a single alphabet for clarity of presentation. However, the same inference pipeline is applicable to all 26 Gujarati alphabet signs included in the dataset.

CONCLUSION

This research successfully developed and validated a robust, high-accuracy framework for static sign recognition of the 26-character Gujarati Sign Language (GSL) alphabet. The primary objective of mitigating data scarcity in regional sign languages was effectively achieved through the application of transfer learning. Comparative evaluation demonstrates that deep CNNs with pre-trained ImageNet weights significantly outperform models trained from scratch. While ResNet50 achieves the highest offline classification accuracy of approximately **99.5%**, MobileNetV2 provides an optimal balance between accuracy (**96.87%**) and computational efficiency, making it more suitable for real-time and mobile deployment. Overall, the proposed framework establishes a reliable and resource-efficient benchmark for GSL recognition, offering strong potential for integration into real-world assistive communication technologies.

Future Work

While the current system provides a strong foundation in static sign recognition, the next phase of research must address the complexities of real-world communication, which is inherently **dynamic and continuous**.

The current system focuses on static GSL recognition; however, real-world sign language communication is dynamic and continuous. Future research must therefore extend the proposed framework to handle temporal variations, large-scale data diversity, and sentence-level interpretation.

- **Dynamic Sign Recognition** is a key extension, where the system must recognise motion-based GSL words and short phrases. This requires integrating temporal sequence models such as LSTM or Transformer architectures with the existing ResNet50 feature extractor. Instead of single images, the model must process sequential feature representations to capture both spatial and temporal characteristics of sign movements.
- **Expansion and standardisation of the GSL dataset** are essential for improving generalisability. Future work should involve collaborative data collection with the GSL Deaf community across different regions to capture signer and dialectal variations. The dataset should be consistently annotated with sign labels, temporal boundaries, and signer metadata to support reproducible research.
- The **development of Continuous Sign Language Recognition (CSLR)** remains the ultimate goal. This involves solving the sign segmentation problem in continuous videos using techniques such as Connectionist Temporal Classification or attention mechanisms, followed by the integration of a Gujarati language model to improve grammatical correctness and contextual accuracy of the translated output.

REFERENCES

1. N. Pigou, S. Dieleman, P.-J. Kindermans, and B. Schrauwen, "Sign language recognition using convolutional neural networks," in Proc. ECCV Workshops, 2014, pp. 572–578.
2. A. Kumar and S. Rani, "Static hand gesture recognition for Indian Sign Language using deep convolutional neural networks," Int. J. Comput. Appl., vol. 182, no. 30, pp. 1–7, 2021.
3. S. Ong and S. Ranganath, "Automatic sign language analysis: A survey and the future beyond lexical meaning," IEEE Trans. Pattern Anal. Mach. Intell., vol. 27, no. 6, pp. 873–891, 2020.
4. M. Mondal, S. Ghosh, and A. Rakshit, "ASL alphabet recognition using transfer learning with VGG16," in Proc. ICCCNT, 2021, pp. 1–6. DOI: 10.1109/ICCCNT51525.2021.9579775
5. G. Thakur and P. Kumar, "Deep residual learning-based Indian Sign Language gesture classification," Multimed. Tools Appl., vol. 82, pp. 21275–21290, 2023.
6. Y. Zhang, H. Huang, and J. Li, "Continuous sign language recognition using hybrid 3D CNN-LSTM networks," IEEE Access, vol. 10, pp. 39812–39824, 2022.
7. M. Tan and Q. V. Le, "EfficientNet: Rethinking model scaling for convolutional neural networks," in Proc. ICML, 2019, pp. 6105–6114. (Used in sign recognition studies 2023–2024.)
8. H. Patel and R. Shah, "Real-time Indian Sign Language digit recognition using CNN," Int. J. Image Process., vol. 16, no. 2, pp. 45–55, 2022.
9. R. Sharma and V. Kumar, "MobileNetV2-based lightweight real-time Indian Sign Language recognition," Comput. Electr. Eng., vol. 110, p. 108803, 2024.
10. M. Alfarraj and A. AlZahrani, "Arabic Sign Language recognition using DenseNet architectures," Sensors, vol. 23, no. 4, p. 2018, 2023.
11. R. Priya and A. Singh, "Challenges in regional sign language dataset construction: A study on Indian regional signs," J. Vis. Commun., vol. 92, p. 103757, 2023.
12. P. Singh, K. Chauhan, and R. Jain, "A multilingual benchmark dataset for South Asian sign languages," IEEE Access, vol. 12, pp. 101245–101260, 2024.
13. S. Sahoo, M. Swain, and S. Mishra, "Dynamic gesture recognition using CNN-BiLSTM hybrid deep networks," Pattern Recognit. Lett., vol. 169, pp. 85–93, 2023.
14. X. Xu, C. Li, and Z. Wu, "Vision transformer-based static sign gesture recognition," Expert Syst. Appl., vol. 237, p. 121643, 2024.
15. M. Rahman, A. Hasan, and T. Ahmed, "Lightweight ASL recognition using MobileNet for edge devices," IEEE Embedded Syst. Lett., vol. 16, no. 1, pp. 45–48, 2024.
16. J. Patel and V. Desai, "YOLOv8-based enhanced Indian Sign Language recognition system for real-time hand detection," IEEE Sensors J., vol. 24, no. 7, pp. 1–10, 2024.