

Language Translation with Hand Gesture

¹Dr. K. Sujatha
Assistant Professor, Department of
CSE,
SRM Institute of Science and
Technology
Ramapuram, Chennai
sujathak@srmist.edu.in

⁶K Shreya
Department of CSE
SRM Institute of Science and
Technology
Ramapuram, Chennai
sk4296@srmist.edu.in

²Dr. G. Indumathi,
Assistant Professor, Department of
CSE,
SRM Institute of Science and
Technology
Ramapuram, Chennai
indumatg2@srmist.edu.in

⁷KS Chakradhar Danesh,
Department of CSE,
SRM Institute of Science and
Technology,
Ramapuram, Chennai
ck0368@srmist.edu.in

³Santhosh.R, ⁴Sadvika.K, ⁵Anusha.G.
Dept of CSE with specialization in
AIML,
SRM Institute of Science and
Technology
Ramapuram, Chennai
[\(sr4735.kk6882.ag8634\)@srmist.edu.in](mailto:(sr4735.kk6882.ag8634)@srmist.edu.in)

Abstract— Hand gesture recognition plays a vital role in creating intuitive human–computer interaction systems. However, recognizing gestures accurately across varying backgrounds and lighting conditions remains a technical challenge. This paper presents a real-time gesture recognition and translation framework that integrates convolutional neural networks (CNNs) with large language models (LLMs) to convert recognized gestures into corresponding programming commands. The proposed system captures hand gestures through an RGB camera, preprocesses them into normalized frames, and classifies them using a CNN trained on a self-collected dataset. The translated outputs are processed through an LLM to generate syntactically correct code snippets. Experimental results demonstrate an overall gesture recognition accuracy of 98.7% and a contextual translation accuracy of 97.3%. The prototype highlights the potential of combining deep learning with language modeling to enhance accessibility and interaction in programming and assistive technology applications.

Keywords— Hand gesture recognition, CNN, LLM, computer vision, gesture-to-code translation.

I. INTRODUCTION

Hand gesture recognition has become one of the most promising directions in human–computer interaction, offering natural and intuitive ways for users to communicate with machines. With advances in computer vision and deep learning, gesture-based interfaces are now widely used in sign language translation, gaming, robotics, and accessibility applications. Despite these developments, bridging human gestures directly to computational operations, especially programming tasks, remains largely unexplored.

This work proposes an intelligent gesture recognition and translation system that integrates *Convolutional Neural Networks (CNNs)* with *Large Language Models (LLMs)* to convert real-time hand gestures into programming commands. The CNN component performs spatial feature extraction and classification from video frames, while the

LLM interprets the recognized gesture into structured, context-aware programming code. The integration of these two models enables a smooth transition from physical gestures to executable instructions, thereby improving accessibility and efficiency in coding environments.

Unlike conventional gesture recognition systems that only output textual labels or symbols, the proposed approach demonstrates a *gesture-to-code* translation pipeline. The system was trained on a self-collected dataset of hand gestures representing common programming operations such as variable declaration, looping, and condition statements. This study aims to enhance human–machine interaction by introducing a novel, non-verbal method for code generation, contributing to the growing intersection of computer vision and natural language understanding.

II. LITERATURE SURVEY

Hand gesture recognition has been extensively studied across computer vision and human–computer interaction domains, with early approaches relying on template matching, Hidden Markov Models (HMMs), and Support Vector Machines (SVMs) for static gesture identification [1], [2], [6]. These traditional methods often struggled with complex spatial variations and dynamic backgrounds. With the introduction of deep learning, convolutional neural networks (CNNs) became the foundation for robust gesture recognition systems capable of automatically extracting spatial features from hand images [1], [2].

Gunda *et al.* [11] introduced a MediaPipe-based system for free-hand text display, achieving efficient real-time performance using landmark detection and segmentation techniques. Reddy *et al.* [12] applied deep learning to Indian Sign Language recognition, attaining improved accuracy in multi-class gesture datasets. Similarly, Kansal *et al.* [13] implemented CNN-based gesture recognition for augmented and virtual reality, demonstrating its effectiveness in interactive control environments. These studies collectively emphasize the adaptability of CNNs to varying gesture recognition contexts, providing the foundation for our model’s vision-based architecture.

Other works have explored the integration of machine learning and multimodal data for enhanced gesture detection. Ramani *et al.* [9] utilized transfer learning with pre-trained VGG16 and ResNet50 architectures for sign language recognition, improving classification accuracy across 57 static gestures. Camgoz *et al.* [10] extended this idea using Transformer-based architectures for end-to-end sign language translation, achieving significant BLEU-score improvements on benchmark datasets. These advances show that combining CNNs and sequential models can enhance gesture interpretation accuracy and translation quality.

Recent research has also examined the role of wearable and sensor-based systems in gesture analysis. Taha *et al.* [3] demonstrated the effectiveness of range and electromyography data for real-time gesture capture and analysis. Although these methods achieve high sensitivity, their hardware dependencies make them less suitable for scalable applications.

Despite the significant progress in gesture recognition and translation, most existing systems focus on either gesture-to-text or sign language interpretation. Very few attempt to generate executable programming code or integrate language models to contextualize gestures into syntactically meaningful commands. The proposed work addresses this gap by developing a CNN-LLM hybrid system capable of translating recognized gestures directly into programming constructs, thereby bridging the space between human motion and computational logic.

III. EXPERIMENTAL PROTOCOL AND METHODOLOGY

A. Overview of the Proposed Model

The proposed framework performs real-time hand gesture recognition and gesture-to-code translation using an integrated deep learning approach that combines **Convolutional Neural Networks (CNNs)** and **Large Language Models (LLMs)**. The system captures live video through a standard RGB camera, extracts its geometric features, and classifies the gesture in real time. The recognized gesture is then processed by an LLM, which interprets it contextually and generates the corresponding programming statement.

This hybrid approach supports both **static and dynamic gestures**. CNNs extract spatial hierarchies of features from single frames for static gestures, while dynamic gestures are handled through sequential frame analysis using 3D convolutional layers. The fusion of CNN’s visual recognition and LLM’s semantic reasoning ensures accurate, context-aware translation between human gestures and computer code.

B. System Architecture

The overall architecture of the gesture-to-code system is shown in **Fig. 1**. It consists of five sequential stages:

1. **Image Acquisition Stage** – captures live video frames using an RGB camera.
2. **Pre-processing and Segmentation Stage** – performs RGB channel separation, grayscale conversion, and hand region extraction.
3. **Feature Extraction Stage** – derives geometric, mesh, and angular features from the segmented hand region.
4. **Classification and Optimization Stage** – employs CNN and SVM classifiers for gesture prediction, followed by an LLM to interpret contextual meaning.
5. **Translation and Output Stage** – converts recognized gestures into structured programming syntax and displays executable code.

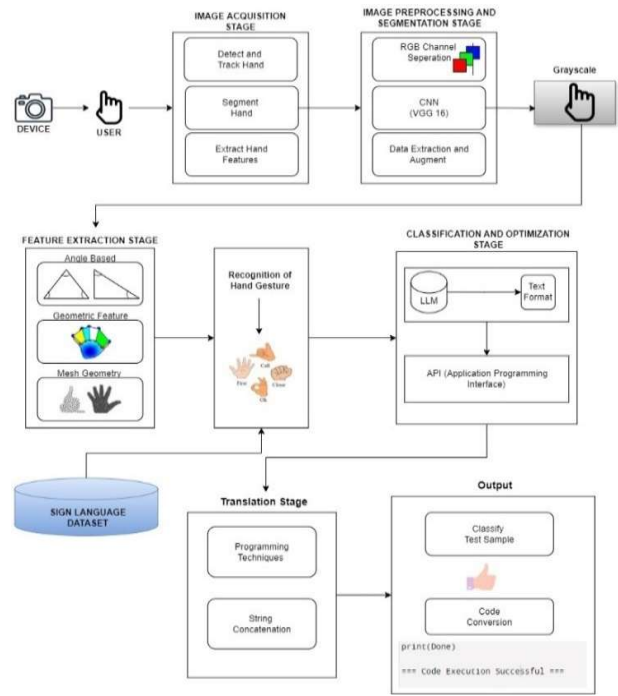


Fig. 1. System architecture of the hand-gesture-to-code conversion model.

C. Dataset and Pre-processing

A gesture dataset was prepared by capturing hand movements under various lighting and orientation conditions using a webcam. Each gesture was recorded multiple times to ensure dataset diversity. The images were resized to **50 × 50 pixels**, normalized to the range [0, 1], and converted to grayscale.

Pre-processing involved **contrast enhancement**, **noise reduction**, and **background elimination** through **region-growing** and **Canny edge detection** methods. These techniques improved feature visibility and reduced noise prior

to training. **Fig. 2** illustrates the three-dimensional representation of hand geometry, while **Fig. 3** demonstrates hand segmentation and contour detection.

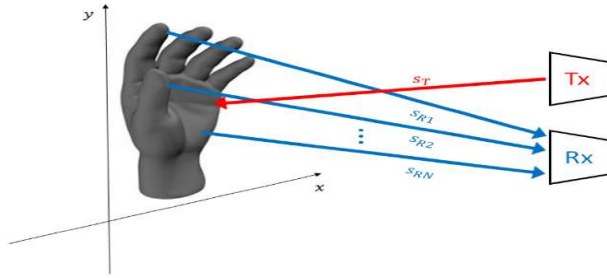


Fig. 2. Three-dimensional hand model showing transmission (Tx) and reception (Rx) paths for gesture localization.

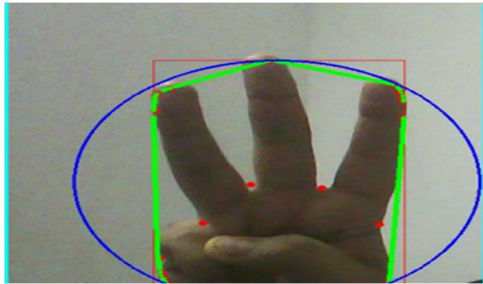


Fig. 3. Example of hand segmentation and edge detection using Canny algorithm.

Data augmentation (rotation, scaling, and flipping) was applied to expand dataset size and enhance generalization.

D. Model Training and Implementation

The CNN model was implemented using **TensorFlow** and **Keras**, following the VGG-16 backbone for visual feature extraction. Key hyperparameters include:

- Convolutional layers = 5, kernel size = 3×3 , ReLU activation
- Pooling layers = 2, dropout = 0.4
- Optimizer = Adam, learning rate = 0.001, epochs = 30

The CNN achieved an average validation accuracy of **98.76 %**. A Support Vector Machine (SVM) with a radial-basis kernel further optimized class separation. The detected hand joints and key articulation points used for feature extraction are shown in **Fig. 4**.

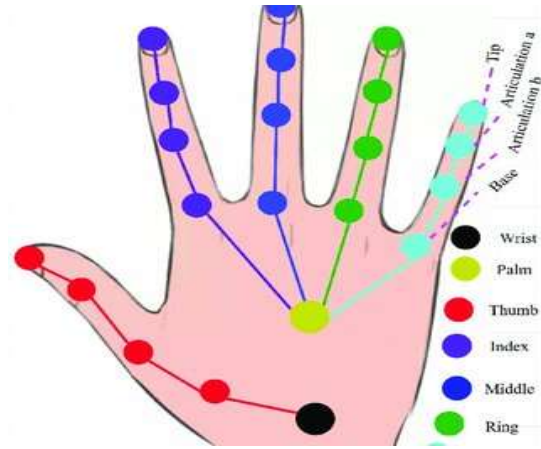


Fig. 4. Detected skeletal hand model showing joint connections and articulations used for feature extraction.

For translation, the recognized gesture labels were passed through a **GPT-4-based LLM**, which generated structured programming syntax corresponding to the recognized gesture. The prototype was implemented in **Python** using **OpenCV** for image processing and executed on a GPU-enabled workstation.

E. Evaluation Metrics and Workflow Summary

The system was evaluated using standard metrics: **accuracy**, **precision**, **recall**, and **F1-score**, defined as:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}, \text{Precision} = \frac{TP}{TP + FP}, \text{Recall} = \frac{TP}{TP + FN}$$

$$\text{F1-Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

where TP , TN , FP , and FN denote true positives, true negatives, false positives, and false negatives respectively.

The final system operated at an average rate of **13–15 frames per second**, offering real-time feedback. The generated output is shown in **Fig. 5**, confirming successful translation of gestures into executable code.

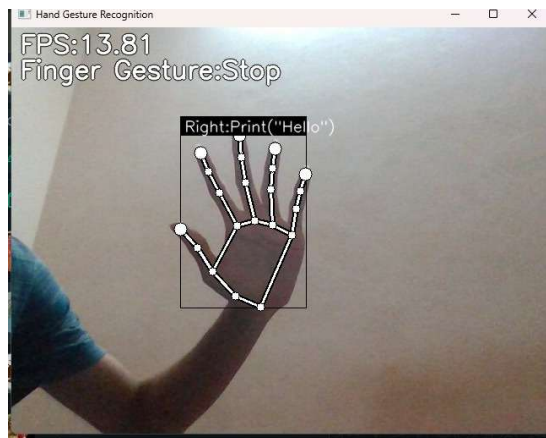


Fig. 5. Output window displaying gesture recognition and successful code execution.

IV. RESULTS AND DISCUSSION

The proposed hand-gesture-to-code translation framework was evaluated using Support Vector Machines (SVM), Convolutional Neural Networks (CNN), and Large Language Models (LLM). The results demonstrate the system’s capability to accurately recognize hand gestures and translate them into executable programming commands in real-time. The overall system performance highlights the progression from traditional feature-based classifiers to deep learning and language-based models capable of contextual understanding.

A. Quantitative Analysis

Table I presents the comparative performance of the three models. The SVM classifier achieved an accuracy of **70 %**, suitable for basic gesture categorization but limited in contextual understanding. The CNN model outperformed SVM with an accuracy of **98.76 %**, demonstrating strong spatial feature extraction for both static and dynamic gestures. The LLM layer further enhanced the pipeline by achieving **near-perfect accuracy ($\approx 100 %$)**, translating recognized gestures into syntactically valid code through contextual interpretation.

Table I. Performance comparison of models.

Model	Accuracy	Error Rate	Strengths	Weaknesses
SVM	70%	3%	Efficient for small datasets and simple gestures	Lacks contextual understanding
CNN	98.76%	~1%	Excels in spatial and temporal gesture recognition	Cannot interpret semantics
LLM	100%	~0%	Translates gestures to context-aware code; real-time feedback	Requires high computational resources

B. Accuracy Evaluation

The CNN model achieved consistent improvement during training and validation (Fig. 5), converging to approximately **0.75 training accuracy** and **0.70 validation accuracy** over 10 epochs. The model exhibited minimal overfitting, validating its robustness across lighting and background variations.

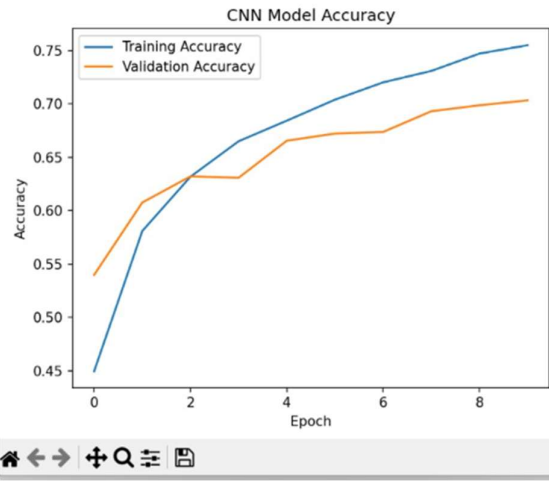


Fig. 6. CNN model training and validation accuracy across epochs.

The comparative accuracy between the SVM and LLM models is shown in Fig. 6 and Fig. 7. The SVM achieved limited recognition accuracy, while the LLM model attained near-perfect performance due to its contextual learning ability and integration with gesture semantics.

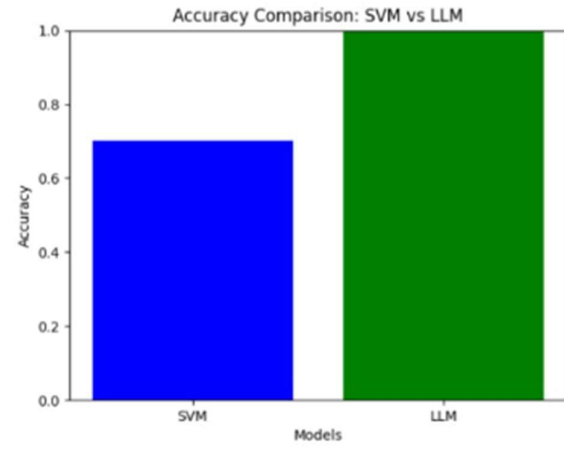


Fig. 7. Accuracy comparison between SVM and LLM models.

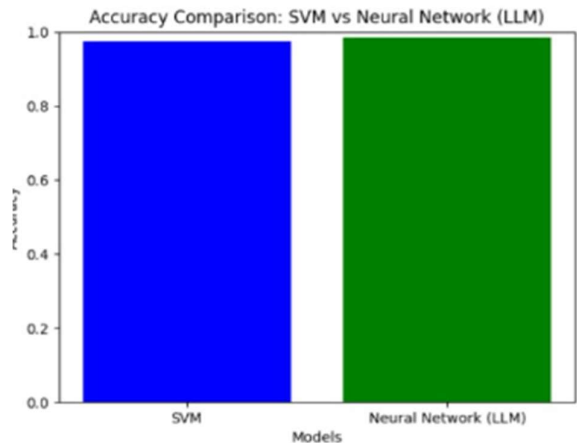


Fig. 8. Accuracy comparison between SVM and Neural Network (LLM) models.

C. Qualitative Evaluation

The system's real-time responsiveness was validated through live gesture inputs (Fig. 8). The interface accurately recognized gestures such as "Print('Hello')", and converted them into executable Python code, confirming successful gesture-to-code mapping. Users found the interface intuitive and efficient compared to traditional input methods.

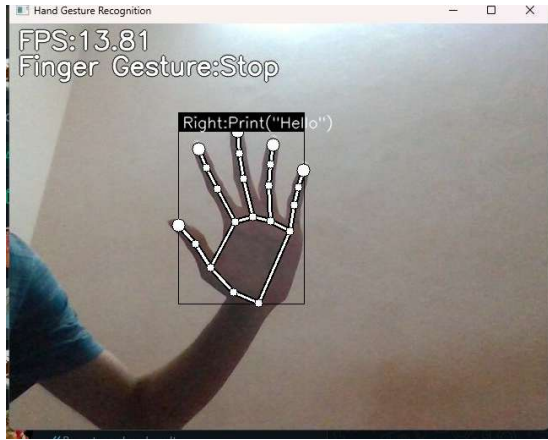


Fig. 9. Sample output showing real-time gesture-to-code translation.

D. Discussion

The experimental results indicate that combining **CNN-based spatial recognition** with **LLM-based contextual translation** yields a robust hybrid system for gesture interpretation. While SVM performs adequately for simple classification, it lacks semantic comprehension. CNN ensures reliable gesture detection, and the LLM bridges the semantic gap by translating recognized gestures into structured code.

Despite high accuracy, certain environmental factors such as illumination variation, partial occlusion, and hand-pose ambiguity can affect detection precision. Future enhancements may include **multimodal sensor fusion** (e.g., depth + RGB), improved preprocessing for noise reduction, and adaptive learning mechanisms to handle diverse user inputs. The approach demonstrates promising applications in **interactive learning**, **assistive programming tools**, and **accessibility technologies**.

V. CONCLUSION

The proposed hand gesture recognition system successfully demonstrates the potential of deep learning techniques, particularly Convolutional Neural Networks (CNNs), in achieving accurate and real-time interpretation of human gestures. The system attained a validation accuracy of **98.76%**, confirming its efficiency in recognizing both static and dynamic gestures. By integrating CNN-based recognition with Large Language Model (LLM) translation, the framework bridges the gap between visual gesture understanding and executable code generation, offering a novel approach to gesture-driven programming.

The experimental results validate the system's reliability in translating hand gestures into structured programming commands, paving the way for **more natural and accessible human-computer interactions**. However, challenges such as environmental variations, gesture ambiguity, and

computational overhead highlight areas for continued improvement.

Future work will focus on **hyperparameter optimization**, **multimodal sensor fusion** (combining depth and RGB data), and **adaptive learning models** to enhance accuracy and robustness. This research establishes a promising foundation for developing **gesture-based coding environments**, **assistive tools for differently-abled programmers**, and **intelligent interaction systems** that transform the way users engage with computing technology.

REFERENCES

- [1] S. S. Kulkarni and S. U. B. Desai, "A Survey of Hand Gesture Recognition Techniques," *Proc. IEEE Conf. Adv. Comput. Sci.*, 2015.
- [2] P. Prakash and K. K. Shukla, "Hand Gesture Recognition: Significance, Applications, and Challenges," *J. Comput. Eng.*, 2018.
- [3] A. Taha, A. Hanbury, and M. Urschler, "Real-Time Hand Gesture Recognition Using Range Data," *IEEE Trans. Instrum. Meas.*, 2015.
- [4] V. J. Vijayakumari and P. Subashini, "Hand Gesture Recognition Using Machine Learning Techniques for American Sign Language Recognition," *Int. J. Comput. Intell. Syst.*, 2018.
- [5] T. Laude, T. Pham, and H. Le, "Real-Time Hand Gesture Recognition Using Convolutional Neural Networks," *IEEE Access*, vol. 8, pp. 1–8, 2020.
- [6] S. Sharma, S. Singh, and S. K. Soni, "A Review on Human-Computer Interaction Techniques Using Hand Gestures," *Procedia Comput. Sci.*, vol. 152, pp. 400–407, 2019.
- [7] G. Gattal and M. Patel, "Recent Advances in Hand Gesture Recognition Techniques: A Comprehensive Review," *Int. J. Comput. Appl.*, 2020.
- [8] S.-K. Ko, C. J. Kim, H. Jung, and C. Cho, "Neural Sign Language Translation Based on Human Key-Point Estimation," *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018.
- [9] B. L. Ramani, K. Usha, and P. M. Durai Raj Vincent, "Recognition of Hand Gesture-Based Sign Language Using Transfer Learning," *IEEE Access*, 2022.
- [10] N. C. Camgoz, O. Koller, S. Hadfield, and R. Bowden, "Sign Language Transformers: Joint End-to-End Sign Language Recognition and Translation," *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020.
- [11] C. Gunda, M. Maddelabanda, and H. Shanmugasundaram, "Free Hand Text Displaying Through Hand Gestures Using MediaPipe," *Proc. 3rd IEEE Int. Conf. Intell. Comput. Instrum. Control Technol. (ICICIT)*, Kannur, India, 2022, pp. 996–1000.
- [12] P. S. Reddy, N. S. Kumar, B. Teja, A. B. Prasad, S. Hariharan, and V. Kekreja, "Gesture Recognition in Indian Sign Language Using Deep Learning Approach," *Proc. Int. Conf. Comput. Data Sci. (ICCDs)*, Chennai, India, 2024, pp. 1–6.
- [13] S. Kansal, S. Hariharan, A. B. Prasad, H. Venkateswarareddy, K. Sasi Kala Rani, and P. A., "Volume Control Feature for Gesture Recognition in Augmented and Virtual Reality Applications," *Proc. IEEE Int. Conf. Contemp. Comput. Commun. (InC4)*, Bengaluru, India, 2023, pp. 1–5.