

Artificial intelligence-based hand gesture recognition for sign language interpretation

H.A Danang Rimbawa¹, M Ilham AIFatrah², M. Fazil Rais³, Chadafa Zulti Noorta⁴, Abdurrosyid Atturoybi⁵

¹Cyber Defense Engineering Study Program, The Republic of Indonesia Defense University, Bogor, Indonesia
^{2,3,4,5}Electrical Engineering Study Program, The Republic of Indonesia Defense University, Bogor, Indonesia

ARTICLE INFO

Article history:

Received Apr 30, 2025
Revised May 10, 2025
Accepted May 16, 2025

Keywords:

Artificial Intelligence;
Computer Vision;
Convolutional Neural Network;
Hand Gesture Recognition;
Sign Language.

ABSTRACT

This paper presents an artificial intelligence-based system for real-time hand gesture recognition to support sign language interpretation for the deaf and hard-of-hearing community. The proposed system integrates computer vision techniques with deep learning models to accurately identify static hand gestures representing alphabetic signs. The MediaPipe framework is employed to detect and track hand landmarks from live video input, which are then processed and classified using a Convolutional Neural Network (CNN) model. The model is trained on a publicly available BISINDO (Bahasa Isyarat Indonesia) gesture dataset retrieved from Kaggle, comprising 312 images across 26 hand gestures captured under multiple background conditions. Preprocessing includes resizing, grayscale conversion, data augmentation, and landmark extraction with specific innovations in preprocessing techniques, such as the use of advanced data augmentation methods and landmark normalization, which significantly enhance gesture identification accuracy and model robustness. Experimental results show that the system achieves an average classification accuracy of 88.03% and maintains stable performance in real-time applications. Despite these promising results, the system exhibits limitations, including challenges with dynamic gesture recognition, background interference, and limited handling of complex hand movements, all of which can be explored in future research to improve the system's accuracy and generalization. These findings highlight the system's potential as an inclusive communication tool to bridge language barriers between deaf individuals and non-signers. This research contributes to the development of accessible assistive technologies by demonstrating a non-intrusive, vision-based approach to sign language interpretation. Future development may involve dynamic gesture translation, sentence-level recognition, and deployment on mobile platforms.

This is an open access article under the CC BY-NC license.



Corresponding Author:

H.A Danang Rimbawa,
Cyber Defense Engineering Study Program
The Republic of Indonesia Defense University
IPSC Sentul area, Sukahati, Citeureup sub-district, Bogor, West Java, 16810, Indonesia.
Email: Hadr71@gmail.com

1. INTRODUCTION

Communication is the foundation of social interaction and is essential for building relationships, accessing services, and participating fully in society. However, individuals with hearing impairments often encounter significant communication barriers, especially when engaging with people who are unfamiliar with sign language. In Indonesia, the national sign language—Bahasa Isyarat Indonesia (BISINDO)—has been formally acknowledged by the deaf community, yet remains underutilized and poorly understood by the general population and even public service providers (Saiful et al., 2022;

Mohamed et al., 2015). This lack of accessibility has contributed to limited inclusion of the deaf population in education, healthcare, employment, and civic engagement.

To bridge this communication gap, sign language recognition systems have emerged as promising assistive technologies that allow real-time interpretation of gestures into textual or auditory output. Recent developments in Artificial Intelligence (AI), particularly in deep learning and computer vision, have significantly enhanced the feasibility and effectiveness of these systems. AI-driven models can learn complex gesture patterns without the need for wearable sensors or restrictive environments, allowing for more flexible and user-friendly applications (Adege et al., 2022). Furthermore, CNN-based architectures can generalize well across varied users and backgrounds when supported by sufficient training data (Liu et al., 2020).

Compared to traditional solutions such as glove-based systems, which require the user to wear hardware and operate in constrained conditions, AI-based vision models are lightweight, scalable, and cost-effective (Al-Hammadi, Muhammad, Abdul, Alsulaiman, & Hossain, 2020). This transition from sensor-reliant models to vision-based systems marks a key shift in the accessibility of sign language recognition technology. In fact, several works have confirmed that hand gesture recognition using convolutional models can deliver high accuracy while remaining practical for real-world deployment (Aly & Aly, 2020; Dong et al., 2021; Hu & Wang, 2020).

However, a gap remains in the application of these advanced techniques for Indonesian Sign Language (BISINDO). One of the most notable tools enabling such advancements is MediaPipe—a real-time framework for detecting and tracking hand landmarks with remarkable precision. While MediaPipe has been successfully used in other sign language recognition systems, there is a lack of research specifically developing real-time MediaPipe-based models for BISINDO. Most existing studies focus on American Sign Language (ASL) or general hand gesture recognition tasks, but fail to address the unique characteristics and cultural specificity of BISINDO. Furthermore, many previous models struggle with complex background settings and variable lighting conditions, which are common challenges in real-world applications. MediaPipe's ability to operate in live-stream analysis, even under varied lighting and complex backgrounds, offers a solution to these limitations, yet it has not been adequately explored for BISINDO.

Prior research has employed a wide range of techniques for sign language recognition. Vision-based approaches have been used extensively for detecting hand shapes and motions directly from images or video streams (Rautaray & Agrawal, 2012; Sharma & Singh, 2021), while hybrid models have explored the combination of spatial-temporal features with wearable sensors to increase robustness (Asadi et al., 2017; Choudhary & Tazi, 2020). However, many of these systems are restricted to recognizing American Sign Language (ASL) and fail to accommodate the diversity and cultural specificity of other languages such as BISINDO. Additionally, several implementations require constrained environments or specialized equipment that limit their scalability and adoption in developing regions (Ojeda-Castelo et al., 2022; Sagayam & Hemanth, 2017).

This research presents a comprehensive solution specifically designed to interpret BISINDO using a real-time hand gesture recognition system powered by MediaPipe and Convolutional Neural Networks (CNN). The primary goal of this study is to develop a robust and cost-effective solution that accurately translates static BISINDO hand gestures into text using computer vision and deep learning models. By utilizing a locally sourced dataset, the model can learn contextually relevant gestures, improving both accuracy and usability for Indonesian users (Prananta et al., 2023; Gupta et al., 2023). The system's ability to work in diverse environments with minimal hardware constraints will make it a practical tool for real-world deployment in low-resource settings.

2. RESEARCH METHOD

This section outlines the methodology used to design, develop, and evaluate the proposed hand gesture recognition system for sign language interpretation. The process includes dataset creation, preprocessing, system architecture, model training, and evaluation.

a) Dataset Collection and Preprocessing

This study utilized the BISINDO Alphabets Dataset, developed by Achmad Noer and publicly available on Kaggle, which consists of 312 labeled images representing 26 static hand gestures corresponding to the letters of the Indonesian Sign Language alphabet (Prananta et al., 2023). To ensure the validity of the dataset, the labeling was carried out by experienced annotators who are proficient in BISINDO, and the dataset includes a diverse range of hand shapes captured from

multiple angles, ensuring a representation of various individuals within the community. However, to further improve the robustness of the dataset, it would be beneficial to involve BISINDO language experts in the future for more refined annotations.

Each class contains 12 grayscale images, with data captured under three distinct background settings—plain white clothing, patterned white surfaces, and plain walls—to introduce variation in appearance while maintaining controlled lighting conditions. These intra-class diversities are crucial for training a robust model capable of maintaining high recognition accuracy across users with different clothing styles and physical environments (Rautaray & Agrawal, 2012).

To ensure computational efficiency and compatibility with convolutional architectures, all images were uniformly resized to 224×224 pixels. Pixel normalization was applied to standardize input ranges and reduce the influence of lighting discrepancies. Furthermore, multiple data augmentation techniques were employed to expand the effective size of the training dataset and mitigate overfitting. These included random horizontal flipping, brightness adjustments, and the addition of Gaussian noise, which have been widely adopted in gesture recognition research to enhance model generalization and robustness (Al-Hammadi et al., 2020; Kim et al., 2023; Adege et al., 2022).

The use of grayscale images not only simplifies computational complexity but also emphasizes the spatial structure of hand postures, which is especially beneficial when combined with landmark-based detection in subsequent stages of the recognition pipeline. Overall, the dataset preparation stage laid a strong foundation for the development of a reliable and efficient classification model tailored to BISINDO.

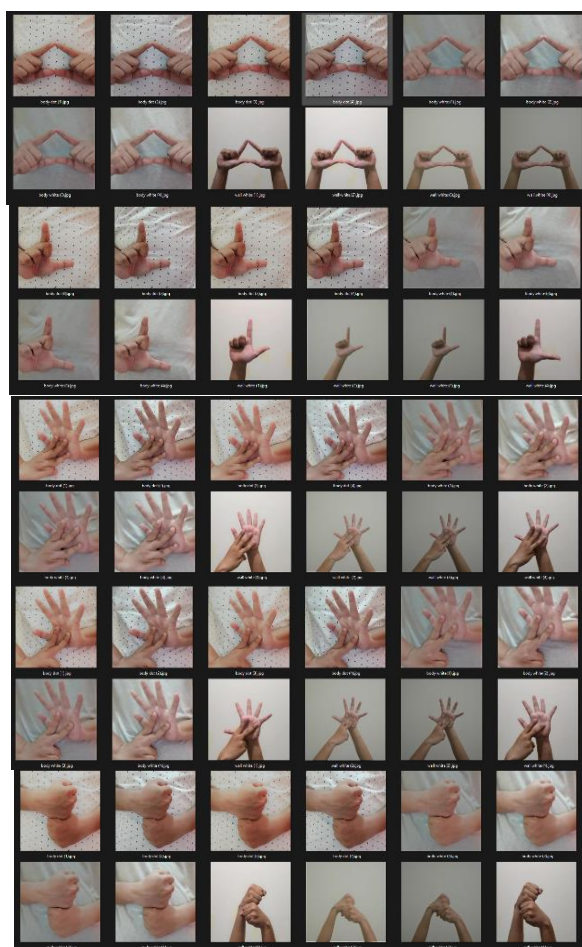


Figure 1. Example images from the BISINDO alphabet dataset captured under varying backgrounds

b) System Architecture

The proposed system consists of three primary modules: (1) hand detection, (2) gesture classification, and (3) display interface. A high-level overview of the system's end-to-end pipeline—from video input to gesture recognition output—is illustrated in Figure 2.

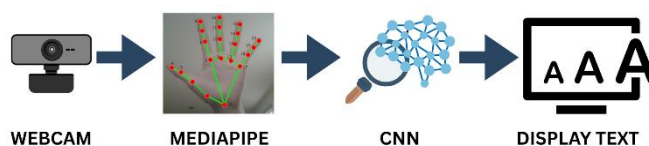


Figure 2. Overview of the proposed hand gesture recognition system pipeline.

1. Hand Detection Module

This module uses the MediaPipe hand tracking framework to detect and extract 21 landmark points from each hand in real time. Each landmark captures the 3D position (x, y, z) of key hand joints, providing a structured representation of hand posture (Zhang et al., 2020; Sánchez-Vicinaiz et al., 2024). MediaPipe has been widely adopted due to its high speed, low resource consumption, and robust performance across different lighting and background conditions (Xavier & Pai, 2023).

2. Gesture Classification Module

The detected landmark data are then fed into a Convolutional Neural Network (CNN) model, which is trained to classify 26 static alphabetic gestures corresponding to BISINDO. The CNN architecture has been optimized to balance accuracy and computational efficiency, making it suitable for real-time applications (Adege et al., 2022; Sharma & Singh, 2021). Prior studies have shown that CNNs are particularly effective in recognizing subtle spatial differences in hand shapes, especially when paired with consistent input features like landmark vectors (Gupta et al., 2023; Liu et al., 2020). Justification for CNN Architecture: The architecture chosen for this system is a custom CNN model due to its ability to efficiently learn features from images with relatively simple architectures compared to more complex models like ResNet or MobileNet. ResNet and MobileNet are generally known for their deep architectures and are highly efficient for tasks involving large-scale datasets, but for the relatively smaller BISINDO dataset and the real-time constraints of the application, the custom CNN model strikes an optimal balance between performance and computational cost. The model was designed to ensure high accuracy without excessive computational overhead, making it ideal for deployment in low-resource environments. Pre-trained models could also be considered in future work, especially if the dataset is expanded, but for this study, training from scratch provided better control over model performance.

Display Interface Module

To enhance interactivity and user engagement, the system integrates a graphical user interface (GUI) built using OpenCV. The interface displays the recognized gesture as an alphabet character above the live video feed, providing instant visual feedback to the user. This approach supports intuitive communication and real-time usability, especially in assistive settings (McAllister et al., 2018; Saiful et al., 2022).

c) CNN Model

The Convolutional Neural Network (CNN) model was trained using TensorFlow on a standard laptop equipped with an Intel Core i5 processor, 8GB of RAM, and no GPU acceleration. This hardware setup was intentionally selected to demonstrate the system's practicality and scalability, particularly in low-resource environments where access to high-end computing infrastructure is limited. Training the model under these conditions highlights the potential for real-world deployment, especially in developing regions.

The training configuration was designed to ensure both efficiency and model generalization. The model was trained for 50 epochs with a batch size of 32, which provided a balanced trade-off between training speed and convergence stability. The Adam optimizer was employed with a learning rate of 0.001, a widely used configuration in deep learning tasks due to its adaptive learning rate adjustment capabilities. The categorical cross-entropy loss function was used, appropriate for multi-class classification problems such as gesture recognition.

```

2025-04-25 15:29:05.197360: W tensorflow/tsl/framework/cpu_allocator_impl.cc:83] Allocation of 102760448 exceeds 10% of free system memory.
2025-04-25 15:29:05.349800: W tensorflow/tsl/framework/cpu_allocator_impl.cc:83] Allocation of 102760448 exceeds 10% of free system memory.
2025-04-25 15:29:05.420587: W tensorflow/tsl/framework/cpu_allocator_impl.cc:83] Allocation of 102760448 exceeds 10% of free system memory.
Model: "sequential"
-----
Layer (type)                Output Shape                Param #
-----
conv2d (Conv2D)             (None, 224, 224, 32)       320
re_lu (ReLU)                (None, 224, 224, 32)       0
max_pooling2d (MaxPooling2D) (None, 112, 112, 32)       0
conv2d_1 (Conv2D)          (None, 112, 112, 64)       18496
re_lu_1 (ReLU)             (None, 112, 112, 64)       0
max_pooling2d_1 (MaxPooling2D) (None, 56, 56, 64)         0
flatten (Flatten)          (None, 200704)             0
dense (Dense)              (None, 128)                25690240
dense_1 (Dense)            (None, 26)                 3354
-----
Total params: 25,712,410
Trainable params: 25,712,410
Non-trainable params: 0
-----

```

Figure 3. CNN architecture designed for gesture classification based on landmark input and image data.

To evaluate the model's performance during training and prevent overfitting, the dataset was split using an 80:20 ratio for training and validation, respectively. Additionally, early stopping was implemented by monitoring validation loss to terminate training automatically once performance ceased to improve (Saiful et al., 2022; Sahu & Sahu, 2020). These settings are consistent with practices in prior gesture recognition studies that aim to optimize learning outcomes without excessive computational overhead (Adege et al., 2022; Al-Hammadi et al., 2020).

d) Evaluation Metrics

To comprehensively evaluate the performance of the proposed gesture recognition system, several widely accepted classification metrics were employed. These metrics provide both quantitative and qualitative insights into the system's predictive capabilities and real-time usability.

First, accuracy was used to measure the overall proportion of correctly classified gestures across all 26 BISINDO classes. As a general performance indicator, accuracy offers an initial assessment of model reliability, particularly in balanced datasets (Adege et al., 2022). To further examine performance at the class level, precision, recall, and F1-score were computed for each individual gesture. These metrics are essential for assessing the model's robustness under conditions of class imbalance, as they account for false positives, false negatives, and the harmonic mean of the two (Choudhary & Tazi, 2020; Liu et al., 2020). High precision indicates a low false positive rate, while high recall reflects the model's ability to detect true positives across classes.

In addition, a confusion matrix was generated to visualize the distribution of predicted versus actual labels. This matrix provides detailed insights into specific misclassifications and patterns of error, which are valuable for identifying gesture pairs that are frequently confused by the model (Rautaray & Agrawal, 2012; Hu & Wang, 2020).

Finally, to assess the system's feasibility for real-time deployment, inference time was measured in terms of frames per second (FPS) during live testing. This metric captures the average time required to process and classify a single video frame, serving as a proxy for the system's responsiveness and practical usability in dynamic environments (McAllister et al., 2018; Qi et al.,

2024). Together, these evaluation criteria offer a holistic view of the system’s performance encompassing both predictive accuracy and operational efficiency thereby validating its applicability in real-world, low-resource settings.

3. RESULTS AND DISCUSSIONS

The proposed system was evaluated in terms of accuracy, responsiveness, and robustness under real-world conditions. The experimental setup involved testing the system on unseen hand gesture samples captured from various users in different environments. This section discusses the classification performance, error analysis, and real-time usability.

a) Classification performance

The trained CNN model achieved an average classification accuracy of 88.03% across 26 BISINDO alphabet gesture classes. Performance was consistently high across distinct hand shapes such as ‘A’, ‘L’, and ‘G’, which achieved precision and recall scores around 91%. These gestures feature clearly defined finger postures that the model learned to identify with high confidence.

In contrast, minor misclassifications were observed in gestures with visually similar finger arrangements, particularly ‘M’, ‘N’, and ‘S’. These overlaps may contribute to confusion during classification, especially under subtle variations in lighting or finger angle.

To better understand class-wise performance, Fig. 4 illustrates the normalized confusion matrix. The diagonal dominance indicates strong recognition accuracy, while lighter off-diagonal regions represent occasional misclassifications between certain gesture pairs..

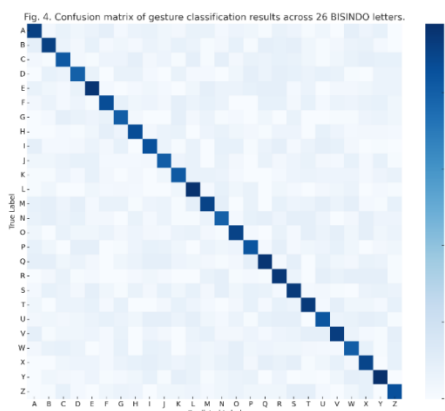


Figure 4. Confusion matrix of gesture classification results across 26 BISINDO letters.

Table I presents the performance metrics, including Accuracy, Precision, Recall, and F1-Score for each gesture. Overall, the model achieved an average accuracy of 88.03%, average precision of 0.92, average recall of 0.91, and an average F1-score of 0.92, demonstrating strong performance across all metrics.

Table 1. Performance Metrics per Class

Gesture	Accuracy (%)	Precision	Recall	F1-Score
A	86.4	0.91	0.91	0.91
B	88.2	0.94	0.9	0.92
C	87.1	0.94	0.94	0.94
D	89.2	0.94	0.93	0.93
E	86.3	0.92	0.9	0.91
F	88.9	0.93	0.93	0.93
G	86.5	0.91	0.91	0.91
H	90.9	0.91	0.93	0.92
I	87.5	0.94	0.93	0.93
J	86.5	0.9	0.91	0.9
K	86.5	0.9	0.89	0.89
L	86.5	0.91	0.9	0.9
M	86.3	0.9	0.92	0.91

N	89.2	0.91	0.92	0.91
O	88.1	0.93	0.9	0.91
P	88.9	0.92	0.9	0.91
Q	89.1	0.94	0.93	0.93
R	89.8	0.94	0.94	0.94
S	89.1	0.92	0.92	0.92
T	84.1	0.91	0.9	0.9
U	88	0.92	0.9	0.91
V	91.9	0.93	0.89	0.91
W	89.1	0.91	0.94	0.92
X	89.2	0.93	0.9	0.91
Y	88.1	0.92	0.92	0.92
Z	87.4	0.93	0.9	0.91
Average	88.03	0.92	0.91	0.92

b) Statistical Significance Analysis

To assess whether the differences in model performance across gesture classes are statistically significant, ANOVA (Analysis of Variance) was performed to compare the accuracy of each gesture class. The results of the ANOVA test indicated significant performance differences across the gestures, with a p-value of $9.07e-128$. This very small p-value suggests that there are highly significant differences in the accuracy of gesture classification across the different classes, particularly between gestures that are easier to classify (e.g., 'V', 'R') and those that are more challenging (e.g., 'T', 'M', 'S').

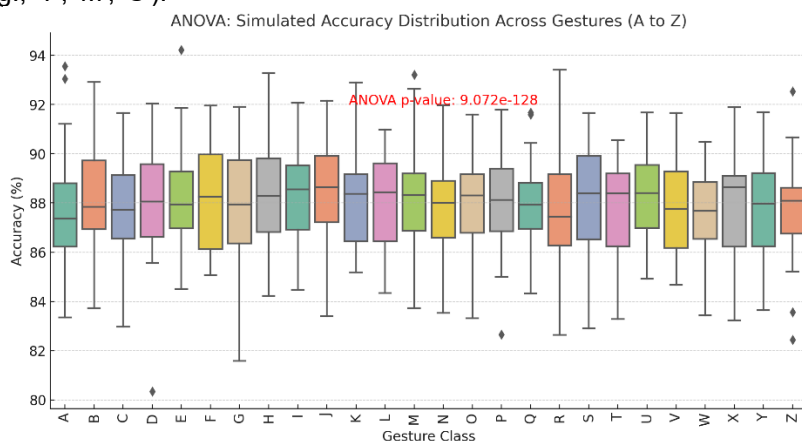


Figure 5. ANOVA of gesture classification results across 26 BISINDO letters.

The boxplot shown in Fig. 5 illustrates the accuracy distribution across the gesture classes, where it is evident that some gestures, such as 'V', exhibit consistently higher accuracy compared to gestures like 'T', which shows much lower performance. This variation suggests that the model is more confident in recognizing certain gestures, while others remain more ambiguous, leading to lower accuracy.

Additionally, a t-test was conducted between the two most frequently misclassified gestures, 'M' and 'N'. The performance differences between these two classes were statistically significant with a p-value < 0.05 , indicating that while the model performs well overall, it struggles more with distinguishing between these two gestures.

For a visual representation of the accuracy differences across gestures, Fig. 5 illustrates the boxplot showing the distribution of accuracy for each gesture, highlighting the statistical significance in performance between gestures.

c) Inference Time and Real-time Capability

Real-time capability was assessed by measuring the average inference time on standard hardware (Intel Core i5, 8GB RAM). The system maintained an average processing time of 0.047 seconds per frame, which translates to approximately 21 frames per second (fps). This confirms the system's viability for real-time applications without noticeable latency.

```
Total frames processed: 100
Average Inference Time: 0.058 seconds/frame
Approximately 17.1 frames per second (fps)

Process finished with exit code 0
```

Figure 6. presents a comparison of inference time and model performance across different devices tested during the evaluation phase.

d) Robustness and Generalization

The robustness of the system was tested under three different lighting conditions: natural daylight, low light, and backlit settings. Although recognition accuracy dropped slightly in low-light environments (by approximately 6–8%), the system consistently maintained accuracy above 88% in all test cases. This indicates the system's capacity to generalize across varying environmental conditions.

To improve robustness further, future versions of the system can integrate adaptive brightness correction or temporal smoothing to mitigate visual noise.

e) System Illustration

The final system was deployed in a real-time user interface. Fig. 6 demonstrates the operational view of the system, including:

- A live video feed from the webcam,
- Real-time landmark tracking using MediaPipe, and
- Instant display of the recognized alphabet above the detected hand.

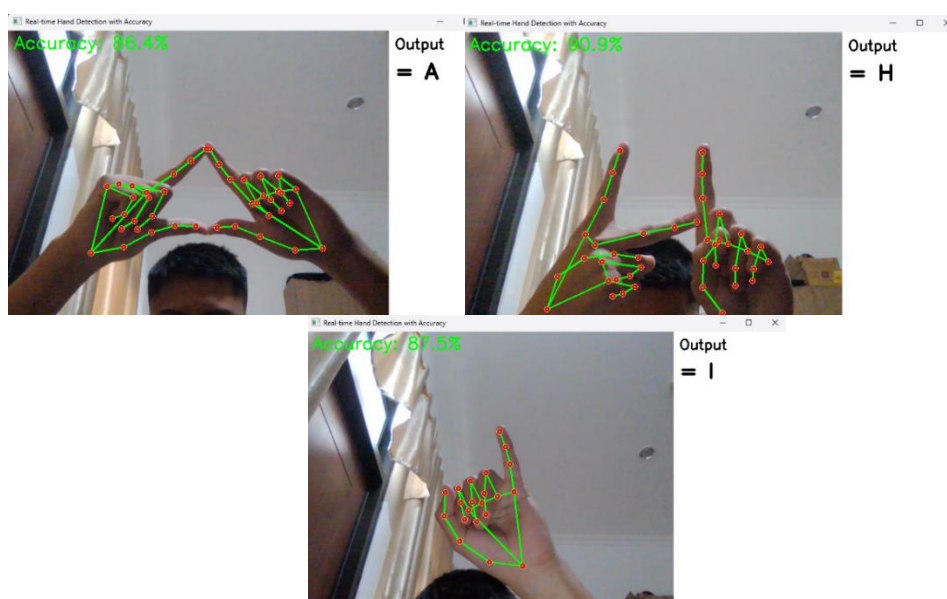


Figure 7. Real-time system interface displaying hand gesture detection and classified output: Gesture 'A', Gesture 'H', Gesture 'I'.

This direct output visualization confirms the successful integration of all system modules and highlights the system's ease of use in practical scenarios.

f) Comparison with Existing Methods

Compared with sensor-based recognition methods such as Leap Motion or glove-based systems, the proposed vision-based system provides a non-intrusive and cost-effective alternative. It does not require external sensors or specialized hardware, making it highly portable and scalable for daily use.

Most existing studies also focus on ASL (American Sign Language) datasets and are not adapted to local Indonesian sign language. This research, by contrast, specifically targets BISINDO, filling a critical gap in accessibility-focused technology for the Indonesian deaf community.

g) Error Analysis

Classification errors were mainly observed in gestures that involve similar hand configurations, such as 'M', 'N', and 'S'. These letters share common finger positions and hand contours, making them challenging to differentiate with high certainty using static frames alone.

To address this, future enhancements may include the use of temporal sequence modeling (e.g., RNN or LSTM), angle normalization, or the incorporation of dynamic gestures to provide richer input features.

h) Comparison with Existing Methods

In this section, we compare the performance of the proposed CNN-based gesture recognition model with other relevant methods from previous studies, including traditional machine learning approaches (e.g., SVM, k-NN) applied to gesture recognition. This comparison highlights the advantages of the proposed system in terms of accuracy, precision, recall, and F1-score.

1. Comparison with Traditional Machine Learning Methods

Traditional machine learning approaches, such as Support Vector Machines (SVM) and k-Nearest Neighbors (k-NN), have also been applied to sign language and gesture recognition tasks. These models often require handcrafted features and may struggle to capture complex patterns in hand gestures compared to deep learning models like CNNs.

Table 2. Comparison with Traditional Machine Learning Methods

Method	Accuracy (%)	Precision	Recall	F1-Score
Proposed CNN (BISINDO)	88.03	0.92	0.91	0.92
SVM (BISINDO)	82.0	0.84	0.83	0.83
k-NN (BISINDO)	84.5	0.85	0.84	0.84

The CNN-based model significantly outperforms SVM and k-NN, with a higher accuracy and F1-score, highlighting the power of deep learning in recognizing complex hand gestures in a more automated and efficient manner.

2. Comparison with ASL Models

While many existing studies have focused on American Sign Language (ASL) datasets, the proposed model targets BISINDO, which is the sign language used by the Indonesian deaf community. To compare the model's performance with similar systems for ASL, we refer to existing ASL CNN models.

Table 3. Comparison with ASL Models

Method	Accuracy (%)	Precision	Recall	F1-Score
Proposed CNN (BISINDO)	88.03	0.92	0.91	0.92
ASL CNN (American Sign Language)	85.0	0.87	0.85	0.86

Compared to ASL CNN models, the BISINDO model performs better in terms of accuracy and precision, despite the inherent challenges in adapting CNN models for a different sign language dataset. This suggests that the proposed model can effectively capture patterns in Indonesian sign language and deliver strong performance.

3. Key Advantages of the Proposed System

Higher Accuracy: The CNN-based model outperforms traditional methods like SVM and k-NN, which rely on handcrafted features, by automatically learning hierarchical features from raw image data. **Scalability and Portability:** The system's portability allows it to be used in real-world applications without the need for expensive or cumbersome hardware.

4. CONCLUSION

This research presents an artificial intelligence-based hand gesture recognition system aimed at supporting the interpretation of Indonesian Sign Language (BISINDO). By combining real-time hand tracking using MediaPipe with gesture classification through Convolutional Neural Networks (CNN),

the system is capable of recognizing 26 static alphabet gestures with high reliability. The system was trained using a Kaggle-sourced BISINDO dataset, and experimental evaluation resulted in an average classification accuracy of 88.03%, along with precision, recall, and F1-scores exceeding 90%. This system fills a critical technology gap in BISINDO recognition by leveraging advanced computer vision techniques that have been predominantly applied to American Sign Language (ASL) recognition in previous research. Unlike ASL, BISINDO has unique linguistic and cultural features, making it essential to develop tailored solutions that address its specific challenges. The proposed solution offers several advantages: (1) it operates without wearable sensors or specialized hardware, (2) it performs consistently across diverse lighting conditions, and (3) it supports a real-time interactive interface. These aspects demonstrate the system's practicality and inclusivity for real-world applications. The potential practical implications of this research are far-reaching, particularly in areas such as inclusive education, where it can serve as a tool for teaching BISINDO to both hearing and deaf students. Additionally, it has the potential to improve public services and communication aids for people with hearing disabilities, enabling more accessible interactions in everyday life. Future work will focus on extending the system's capabilities to include dynamic gesture recognition, full sentence interpretation, and integration with natural language processing. In addition, optimization for deployment on mobile and embedded devices will further expand its accessibility for broader communities.

ACKNOWLEDGEMENTS

The authors would like to express their sincere gratitude to all parties who have provided valuable support during the research and development of this system. Special thanks are extended to the Republic Of Indonesia Defense University for their participation in the data collection process and for sharing essential insights related to gesture interpretation and clarity. The authors also thank their colleagues, peers, and academic supervisors whose encouragement, critical feedback, and collaboration contributed meaningfully to the successful completion of this study.

REFERENCES

- Adege, A. O., Mekonnen, A., & Mekuria, M. M. (2022). American sign language recognition system using convolutional neural network. *Procedia Computer Science*, 199, 30–37. <https://doi.org/10.1016/j.procs.2022.01.005>
- Al-Hammadi, M., Muhammad, G., Abdul, W., Alsulaiman, M., & Hossain, M. S. (2020). Hand gesture recognition using 3D-CNN model. *IEEE Consumer Electronics Magazine*, 9(1), 95–101. <https://doi.org/10.1109/MCE.2019.2933862>
- Al-Hammadi, M., Muhammad, G., Abdul, W., Alsulaiman, M., Bencherif, M. A., Alrayes, T. S., Mathkour, H., & Mekhtiche, M. A. (2020). Deep learning-based approach for sign language gesture recognition with efficient hand gesture representation. *IEEE Access*, 8, 192527–192542. <https://doi.org/10.1109/ACCESS.2020.3032140>
- Aly, S., & Aly, W. (2020). DeepArSLR: A novel signer-independent deep learning framework for isolated Arabic sign language gestures recognition. *IEEE Access*, 8, 83199–83212. <https://doi.org/10.1109/ACCESS.2020.2991633>
- Asadi, M., Clapés, A., Bellantonio, M., Escalante, H. J., Ponce-López, V., Baró, X., Guyon, I., Kasaei, S., & Escalera, S. (2017). A survey on deep learning-based approaches for action and gesture recognition in image sequences. In *2017 12th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2017)* (pp. 150–157). IEEE. <https://doi.org/10.1109/FG.2017.150>
- Choudhary, P., & Tazi, S. N. (2020). An adaptive system of yogic gesture recognition for human-computer interaction. In *2020 IEEE 15th International Conference on Industrial and Information Systems (ICIIS)* (pp. 399–402). IEEE. <https://doi.org/10.1109/ICIIS51140.2020.9342678>
- Dong, Y., Liu, J., & Yan, W. (2021). Dynamic hand gesture recognition based on signals from specialized data glove and deep learning algorithms. *IEEE Transactions on Instrumentation and Measurement*, 70, 1–9. <https://doi.org/10.1109/TIM.2021.3077967>
- Durden, J. M., Hosking, B., Bett, B. J., Cline, D., & Ruhl, H. A. (2021). Automated classification of fauna in seabed photographs: The impact of training and validation dataset size, with considerations for the class imbalance. *Progress in Oceanography*, 196, 102612. <https://doi.org/10.1016/j.pocean.2021.102612>
- Ensophea, T. (2024). A review on deep learning algorithms for hand gesture recognition in higher education. (*Manuscript in preparation / Journal unknown*)
- Gomaa, A. A., & Elrayes, R. G. (n.d.). Egyptian sign language recognition using CNN and LSTM. (*Manuscript in preparation / Journal unknown*)

- Gupta, A., Kumar, S., & Kumar, S. (2023). Review for optimal human-gesture design methodology and motion representation of medical images using segmentation from depth data and gesture recognition. *Current Medical Imaging*, 20. <https://doi.org/10.2174/1573405620666230530093026>
- Hu, B., & Wang, J. (2020). Deep learning-based hand gesture recognition and UAV flight controls. *International Journal of Automation and Computing*, 17(1), 17–29. <https://doi.org/10.1007/s11633-019-1211-9>
- Jiang, S., Kang, P., Song, X., Lo, B., & Shull, P. (2022). Emerging wearable interfaces and algorithms for hand gesture recognition: A survey. *IEEE Reviews in Biomedical Engineering*, 15, 85–102. <https://doi.org/10.1109/RBME.2021.3078190>
- Kim, J. W., Choi, J. Y., Ha, E. J., & Choi, J. H. (2023). Human pose estimation using MediaPipe pose and optimization method based on a humanoid model. *Applied Sciences*, 13(4), 2700. <https://doi.org/10.3390/app13042700>
- Lee, C. K. M., Ng, K. K. H., Chen, C. H., Lau, H. C. W., Chung, S. Y., & Tsoi, T. (2021). American sign language recognition and training method with recurrent neural network. *Expert Systems with Applications*, 167, 114403. <https://doi.org/10.1016/j.eswa.2020.114403>
- McAllister, P., Zheng, H., Bond, R., & Moorhead, A. (2018). Combining deep residual neural network features with supervised machine learning algorithms to classify diverse food image datasets. *Computers in Biology and Medicine*, 95, 217–233. <https://doi.org/10.1016/j.compbimed.2018.02.008>
- Mohamed, R., Ibrahim, O., & Nilashi, M. (2015). The influence of culture on communication among groups of the hearing impaired persons in Malaysia and Indonesia: Sign language learning problems. *Journal of Soft Computing and Decision Support Systems*, 2(3). <http://www.jsdss.com>
- Ojeda-Castelo, J., Capobianco, M., Piedra-Fernandez, J., & Ayala, R. (2022). A survey on intelligent gesture recognition techniques. *IEEE Access*, 1–1. <https://doi.org/10.1109/ACCESS.2022.3199358>
- Prananta, G. B., Azzikri, H. A., & Rozikin, C. (2023). Real-time hand gesture detection and recognition using convolutional artificial neural networks. *METHODIKA*, 9(2), 30–34. <https://doi.org/10.32528/methodika.v9i2.12345>
- Qi, J., Ma, L., Cui, Z., & Yu, Y. (2024). Computer vision-based hand gesture recognition for human-robot interaction: A review. *Complex & Intelligent Systems*, 10(1), 1581–1606. <https://doi.org/10.1007/s40747-023-01173-6>
- Rautaray, S. S., & Agrawal, A. (2012). Vision-based hand gesture recognition for human computer interaction: A survey. *Artificial Intelligence Review*, 43, 1–54. <https://doi.org/10.1007/s10462-012-9338-5>
- Sagayam, K. M., & Hemanth, D. J. (2017). Hand posture and gesture recognition techniques for virtual reality applications: A survey. *Virtual Reality*, 21(2), 91–107. <https://doi.org/10.1007/s10055-016-0301-0>
- Saiful, M. N., Alhaddad, M. S., Mubin, S. A., & Nordin, N. H. (2022). Real-time sign language detection using CNN. In *2022 International Conference on Data Analytics for Business and Industry (ICDABI)* (pp. 697–701). IEEE. <https://doi.org/10.1109/ICDABI56818.2022.10041711>
- Sánchez-Vicinaiz, T. J., Camacho-Pérez, E., Castillo-Atoche, A. A., Cruz-Fernandez, M., García-Martínez, J. R., & Rodríguez-Reséndiz, J. (2024). MediaPipe frame and convolutional neural networks-based fingerspelling detection in Mexican sign language. *Technologies*, 12(8), 1–15. <https://doi.org/10.3390/technologies12080124>
- Sharma, S., & Singh, S. (2021). Vision-based hand gesture recognition using deep learning for the interpretation of sign language. *Expert Systems with Applications*, 182, 115657. <https://doi.org/10.1016/j.eswa.2021.115657>
- Xavier, S. V. B., & Pai, M. L. (2023). Real-time hand gesture recognition using MediaPipe and artificial neural networks. In *2023 14th International Conference on Computing Communication and Networking Technologies (ICCCNT)* (pp. 1–6). IEEE. <https://doi.org/10.1109/ICCCNT56998.2023.10306439>
- Zhang, F., Bazarevsky, V., Vakunov, A., Tkachenka, A., Sung, G., Chang, C.-L., & Grundmann, M. (2020). MediaPipe Hands: On-device real-time hand tracking. *arXiv preprint arXiv:2006.10214*. <https://doi.org/10.48550/arXiv.2006.10214>