

# Hand Gesture Recognition: Enhancing Accuracy and Precision with Deep Learning

1<sup>st</sup> T. G. Ramnadh Babu

*Asst. Professor, Dept of CSE*  
Narasaraopeta Engineering College  
Narasaraopet, Palnadu, AP, India  
baburamnadh@gmail.com

2<sup>nd</sup> Dodda Venkata Reddy

*Asst. Professor, Dept of CSE*  
Narasaraopeta Engineering College  
Narasaraopet, Palnadu, AP, India  
doddavenkatareddy@gmail.com

3<sup>rd</sup> Dr. S.V.N. Sreenivasu

*Professor, Dept of CSE*  
Narasaraopeta Engineering College  
Narasaraopet, Palnadu, AP, India  
drsvnsrinivasu@gmail.com

4<sup>th</sup> K. Srinivasa Kalyan Ram

*Dept of CSE*  
Narasaraopeta Engineering College  
Narasaraopet, Palnadu, AP, India  
ksrinivasakalyanram@gmail.com

5<sup>th</sup> Venkata Aravind Tunuguntla

*Dept of CSE*  
Narasaraopeta Engineering College  
Narasaraopet, Palnadu, AP, India  
venkataaravind2021@gmail.com

6<sup>th</sup> Shaik Hazare Khaja Mohiddin

*Dept of CSE*  
Narasaraopeta Engineering College  
Narasaraopet, Palnadu, AP, India  
mohiddins511@gmail.com

**Abstract**—Accurate, real-time recognition of hand gestures in dynamic environments remains challenging in human-computer interaction. This paper presents a hybrid deep learning model combining Convolutional Neural Networks (CNN) with Recurrent Neural Networks (RNN) using Long Short-Term Memory (LSTM) layers to capture both spatial and temporal information for dynamic hand gesture recognition. Trained on a dataset of six gestures—scroll-left, scroll-right, scroll-up, scroll-down, zoom-in, and zoom-out—the model achieves an accuracy of 95.59%, with an F1 score of 0.94 and an AUC-ROC of 0.95, indicating significant improvement over traditional models and practical viability in real-world applications. Key topics include data preprocessing, model architecture, hardware and software configurations, and performance comparisons with benchmarks. The paper concludes with discussions on limitations and future research directions to enhance the model’s adaptability and efficiency.

**Index Terms**—Human-computer interaction, hand gesture recognition, dynamic hand gestures, video-based gesture recognition, machine learning, deep learning, convolutional neural networks (CNN), recurrent neural networks (RNN), long short-term memory (LSTM).

## I. INTRODUCTION

With information technology progressing rapidly, HCI has grown manifold with an increased demand for more natural modes of interaction, such as gesture-based communication. Hand gestures are a medium of non-verbal communication used in numerous applications like smart homes, gaming, learning, and control of automotive systems. In the early days, gesture recognition systems relied on external sensors like gloves, which could not be put to widespread use. Indeed, modern vision-based systems, utilizing video cameras and deep learning algorithms, especially convolutional neural networks, have been highly successful in the area of gesture recognition. Despite these advances, several challenges remain. Current systems struggle to process real-time data and recognize patterns of complex and rare gestures in dynamic, real-world environments. Additional issues include latency

and the simultaneous recognition of multiple gestures, which limit real-world applications. Moreover, most systems perform poorly when generalizing to unseen gestures and under different conditions or noisy states [1], [3], [9], [10], [11], [13]. Recent works are beginning to approach gesture recognition in relatively unconstrained, “in the wild” settings. Using minimal external equipment in such scenarios can help reduce noise and motion-tracking errors, but poses significant challenges for research due to the lack of large-scale, publicly available 3D skeletal hand datasets. Most works focus on hand-centric rather than full-body gestures, and models trained in ideal conditions often face performance gaps in practical applications with background noise or variable lighting that affect recognition accuracy [4], [2], [7]. In this paper, we recommend a hybrid architecture for deep learning that combines CNNs with Recurrent Neural Networks to improve dynamic hand gesture recognition in natural environments. Using Inception-v3 for feature extraction and LSTM layers for temporal processing, the model’s average accuracy reached 95.59%. The rest of the paper is organized as follows: Section 2 covers related work, Section 3 introduces the proposed method, outlining data preprocessing and model design, Section 4 presents results and discussion, and Section 5 covers the main findings, limitations, and future research directions.

## II. RELATED WORK

The latest advances in CNNs have merged feature extraction and classification into unified architectures, thereby improving recognition tasks for images and videos [3], [5]. Early CNNs handled simpler tasks, like recognizing handwritten digits in the 1990s, but struggled with complex image and video classification until AlexNet achieved notable success on the ImageNet dataset [13]. This success led to deeper architectures such as VGGNet, GoogleNet, and ResNet, which provided greater classification performance but with higher computational costs, making them impractical for resource-limited

Work (Year)	Dataset	Size	Deep Learning Method	Environment	Accuracy
Skeleton-Based Dynamic Hand Gesture Recognition Using PB-GRU-RNN (2020)	SHREC'17 dataset	2,800 gestures	Part-Based GRU-RNN (PB-GRU-RNN)	Lab-like	90.50%
Dynamic Hand Gesture Recognition Using Multi-Branch Attention (2023)	Three dynamic skeleton-based gesture datasets	Not provided	Multi-Branch Attention Based Graph Neural Networks	Real	91.40%
3D Convolutional Neural Networks for Dynamic Gesture Recognition (2017)	20BN-jester	148,092 videos	3D CNN	Real	93.60%
Multi-modal Hand Gesture Recognition Using RGB-D Data (2019)	MSR Gesture 3D	12,000 gestures	3D CNN + LSTM	Lab-like	95.20%
Dynamic Hand Gesture Recognition Using 3D CNN and LSTM (2021)	EgoGesture	24,161 videos	3D CNN + LSTM	Lab-like	92.10%

TABLE I  
COMPARISON OF VARIOUS WORKS ON DYNAMIC HAND GESTURE RECOGNITION

devices [1], [4]. To address this, lightweight CNN models have been developed for effective real-world applications [6]. Most video classification techniques use 2D CNNs to capture spatial features, while hybrid techniques, such as 3D CNN and CNN-RNN, balance computational efficiency with spatial-temporal accuracy [2], [7]. Techniques like Temporal Segment Networks (TSN) and Temporal Shift Modules (TSM) further advanced real-time gesture recognition, with TSM showing higher accuracy in some cases [8]. Much of the related work focuses on static and dynamic hand gesture recognition, with dynamic recognition being more challenging due to the real-time capture of hand shape, movement, and direction [11]. Research suggests that combining CNNs with LSTMs enhances dynamic gesture recognition by leveraging spatial and temporal features [10]. Two-stream CNNs, which process features across spatial and temporal domains, have shown higher accuracy than single-stream models, with reduced computational costs [9]. Some methods apply 3D CNNs directly to video sequences for spatiotemporal feature capture, although performance may degrade under complex scenarios involving background changes or lighting variations. Robustness is improved when combining 3D CNNs and LSTMs, which capture both short-term and long-term spatiotemporal features [14]. Dataset quality and availability, as shown in Table 1, also play crucial roles; datasets captured in realistic settings, such as ChaLearn LAP IsoGD and IPN Hand, may lower accuracy due to environmental complexities, while simpler backgrounds often yield higher accuracy [12], [15].

### III. METHOD

#### A. HYBRID DEEP LEARNING ARCHITECTURE

Videos capture gesture movements more effectively than static images due to their sequential frame arrangement, which is crucial for recognizing motion [5]. To leverage this, we employ a hybrid architecture combining Convolutional Neural Networks (CNN) with Recurrent Neural Networks (RNN) through an LSTM layer [13]. Inception-v3 extracts features from individual frames, which are then used by the RNN for classification [3]. Representations of video files using 2-D spatiotemporal feature maps obtained by CNNs have shown good performance for various video classification tasks

[7]. This hybrid architecture, CNN-RNN, uses CNNs for feature extraction in space and RNNs for temporal sequence processing. We utilized the Inception-v3 model, which has gained widespread acceptance due to its efficiency, for feature extraction from every frame. We removed the last classification layers of the model so we could make use of the output features as inputs to the RNN. Every frame output, of size 2048 features, is fed into the RNN as a temporal sequence of feature maps with a dimensionality of (None, 10, 2048), where "None" indicates variable batch sizes. Most CNNs contain a convolution layer, pooling layer, and fully connected or dense layer. The convolution layer applies a kernel to the 2D grid of an image, hence extracting features from it. Multi-parallel pipelines inside an inception block map the features from several scales in the Inception-v3 model. The great reduction in computational cost enables the model to catch both big and small patterns while maintaining essential information. Output from the last inception block feeds a global average pooling layer reducing the dimensionality before feeding into the RNN. The LSTM layer in the network is an RNN component, capturing the temporal context present in the frame sequences. The LSTMs are one of the advanced forms of RNNs, which solve the issues of exploding or vanishing gradients. Thus, they are proper for learning long-term dependencies of sequence data. To prevent overfitting, the model utilizes dropout layers within both LSTM and fully connected layers, with final classification done through a softmax activation function. This hybrid approach contrasts with other techniques, such as two-stream networks that process spatial and temporal data in parallel, or 3D CNNs that directly analyze video sequences. Our method leverages CNNs and RNNs sequentially to efficiently capture both spatial and temporal features. As illustrated in Figure 1, the architecture effectively integrates these components. The following pseudo-code outlines the core steps involved in implementing the CNN-LSTM architecture for dynamic gesture recognition, covering preprocessing, feature extraction, sequence modeling, and classification.

#### B. EVALUATION

Model performance is generally assessed by accuracy, the ratio of true positive and true negative instances to the total population. Additionally, a confusion matrix provides deeper

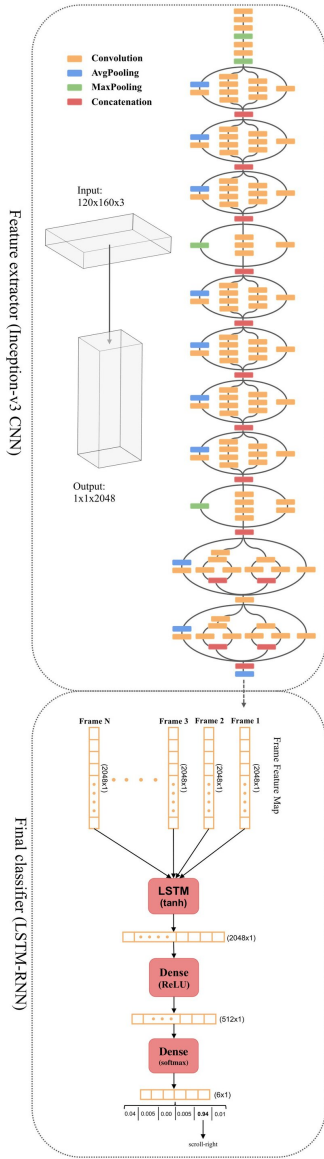


Fig. 1. CNN-RNN Hybrid architecture[13]

insights into precision and sensitivity (recall). Precision is the ratio of true positives to positive predictions, while sensitivity refers to the proportion of true positives out of actual positives [8].

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN} \quad (1)$$

$$\text{F1-score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (2)$$

$$\text{AUC} = \int_0^1 \text{TPR(FPR)} d(\text{FPR}) \quad (3)$$

We treat all six gesture classes equally and adopt a k-fold cross-validation approach for robust performance evaluation. Other RNN architectures, such as GRUs, will be tested for comparison, along with other feature extraction models like

#### Algorithm 1 Dynamic Gesture Recognition with CNN-LSTM

```

0: Input: Video dataset of gesture sequences, pre-trained CNN model (Inception-v3)
0: Output: Predicted gesture class
0: procedure PREPROCESSDATA
0:   for each video in dataset do
0:     Extract frames
0:     Resize frames to (120, 160, 3)
0:     Normalize pixel values
0:   end for
0:   return preprocessed frames
0: end procedure
0: procedure TRAINMODEL
0:   for each frame sequence in preprocessed data do
0:     for each frame in sequence do
0:       Pass frame through CNN (Inception-v3) to extract features
0:       Append extracted features to feature sequence
0:     end for
0:     Pass feature sequence through LSTM layers
0:     Classify sequence with softmax layer
0:   end for
0: end procedure
0: procedure EVALUATEMODEL
0:   Calculate Accuracy, F1 Score, and AUC on test data
0:   Analyze Confusion Matrix for misclassification patterns
0: end procedure
0: Output: Classification accuracy, F1 Score, AUC, and Confusion Matrix =0

```

VGG16, VGG19, ResNet50, and Xception [4]. This allows us to evaluate the robustness of our approach across different model settings.

#### C. HAND GESTURE DATASET

The Depth\_Camera\_Dataset is used for training and evaluation, consisting of six hand gestures with 662 sequences each. Each sequence is a 40-frame set depicting one of the following actions: scroll-left, scroll-right, scroll-up, scroll-down, zoom-in, and zoom-out. This dataset includes various individuals performing the same gestures in realistic settings, capturing natural movements [14]. Figure 2, the dataset structure and the process for extracting frames are depicted.

We use a total of 3000 sequences—500 sequences per class—to manage computational constraints and extract every fourth frame from each sequence [13]. The data is divided into 70% for training, 20% for validation, and 10% for test sets, ensuring class balance and random distribution [5]. K-fold cross-validation is performed with the dataset divided into 5 stratified folds [1].

#### D. DATA PRE-PROCESSING

Data preprocessing utilizes Python libraries such as NumPy, Pandas, and TensorFlow. Gesture frames are stored in a 4-dimensional array within a Pandas DataFrame. RGB values

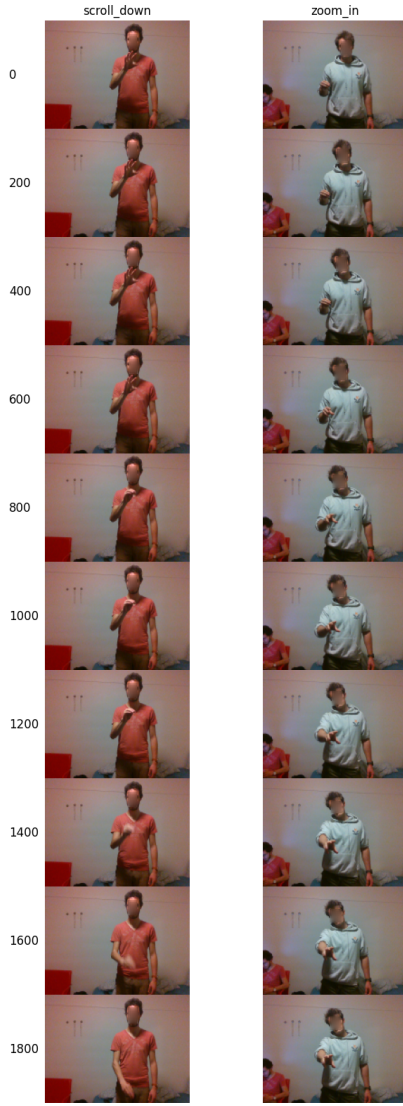


Fig. 2. Example of RGB dynamic hand gestures: scroll-down and zoom-in.

are encoded, resulting in a final input shape of (10, 120, 160, 3) for ten frames per gesture [5]. To balance accuracy with processing time, every fourth frame is used [13]. The dataset is shuffled before splitting to maintain uniform distribution across gestures. K-fold cross-validation is conducted for robust evaluation, averaging results across different folds [1]. Instances are provided unbatched, as final model fitting did not benefit from batching [3].

#### IV. RESULTS AND DISCUSSION

This end-to-end machine learning project was developed and run on Colab, using robust GPU resources for the deep learning workloads. The model was fitted on an NVIDIA A100-SXM GPU via TensorFlow v2.11.0. Training of the model was done at a 0.001 learning rate by categorical cross-entropy as a loss function and Adam as an optimizer. Early stopping was used to prevent overfitting. Besides accuracy, the performance of the model is verified in terms of precision,

recall, F1-score, and AUC-ROC. These metrics help explain better with respect to the balance of classes and their respective performance better than with this single measure. The model attained average precision of 0.93 and an average F1-score of 0.94 across all five folds, which would reflect a fine balance between precision and recall, suggesting that the model can do a good job in terms of distinguishing between the six classes of gestures. Average AUC-ROC scores, which directly reflect the ability of the classifier to distinguish between positive and negative classes at a wide range of thresholds, average at 0.95, with high discriminatory power.

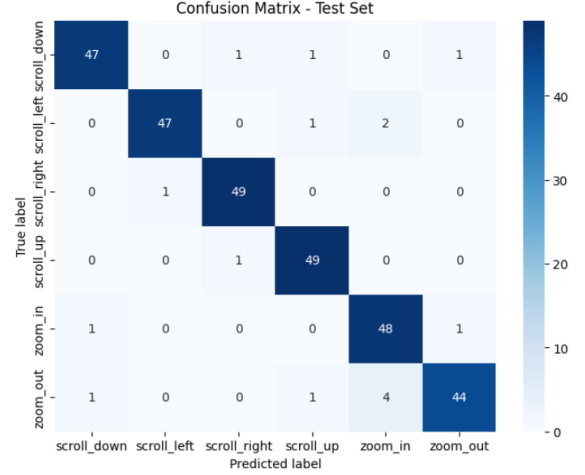


Fig. 3. Confusion matrix of test-data fold 1

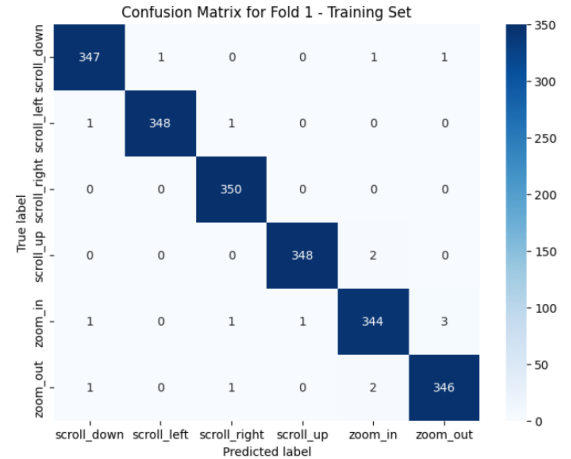


Fig. 4. Confusion matrix of training-data fold 1

##### A. Confusion Matrix and Analysis

The confusion matrix in Figure 3 shows that the model performed best on the "scroll-right" gesture, achieving a True Positive Rate (TPR) of 98%. The lowest TPR was observed for the "zoom-out" gesture, at 88%. Misclassifications occurred primarily between "scroll-left" and "zoom-out," with approximately 4% of "scroll-left" instances misclassified as "zoom-

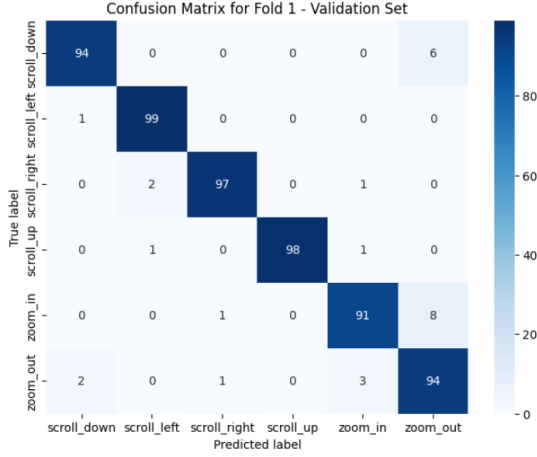


Fig. 5. Confusion matrix of validation-data fold 1

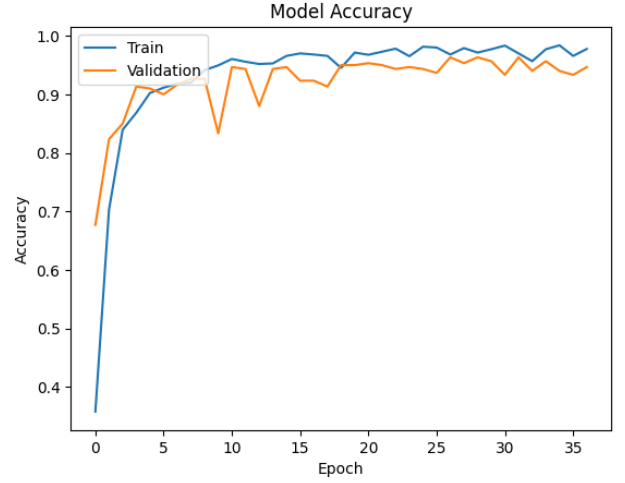


Fig. 7. Accuracy during training

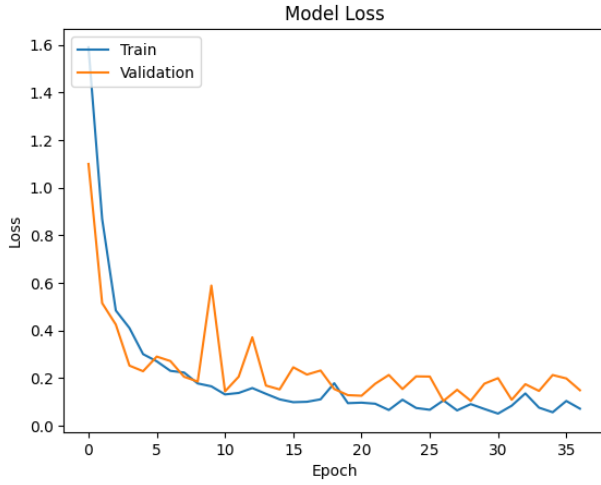


Fig. 6. Loss During Training.

out.” Additionally, 8% of ”zoom-in” instances were misclassified as ”zoom-out.” These errors may be attributed to the subtle visual similarities and overlapping movements between certain gestures, particularly between ”zoom-in” and ”zoom-out,” as these two gestures involve similar directional hand motions, which can lead to ambiguity in model predictions. The same was true in the training and validation datasets presented by Figures 4 and 5, respectively.

### B. Inception-v3 and LSTM Model Performance

The inception-v3 architecture does perform with significant success in pulling features at multiple scales, thus gaining popularity rather than the simple architecture of VGG16 and ResNet50. This is evident from the higher precision and recall scores. In order to maintain care of long-term dependencies in gesture sequences, LSTM layers have been incorporated instead of GRU. Though GRUs are computationally less expensive, LSTMs performed well in handling the temporal sequence of frames, which is a precursor for dynamic gestures recognition.

### C. Training and Validation Performance

Figures 3, 4, and 5 show the confusion matrices for training, validation, and testing, respectively. Overall, the model achieved 100% training accuracy, while maintaining an average validation accuracy of around 95.19% across the folds, indicating good generalization without overfitting. Additionally, to prevent further training once validation loss stabilized, early stopping was employed. Figures 6 and 7 illustrate the model’s accuracy and loss during training.

### D. Comparative Analysis

From the comparative analysis, it is noted that the Inception-LSTM model achieves superior performance for dynamic hand gesture recognition, with an average test accuracy of 95.59%. This model leverages the Inception module’s multi-scale spatial feature extraction and LSTM’s temporal analysis capabilities, making it particularly effective for recognizing gestures with subtle differences, such as ”scroll-up” versus ”scroll-right.” Although the Inception-GRU model provides a similar accuracy with lower computational costs due to the efficiency of GRU, the Xception-LSTM and Xception-GRU models display lower accuracies, around 89-91%, which may be due to their simpler architectures. Consequently, the Inception-LSTM model strikes the best balance between accuracy and feature capture, making it especially suitable for detailed and dynamic gesture recognition.

Model	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Average
Train Accuracy (Inception-LSTM)	100.00%	99.64%	99.94%	99.88%	100.00%	99.89%
Validation Accuracy (Inception-LSTM)	0.01%	0.03%	0.01%	0.01%	0.00%	99.98%
Test Accuracy (Inception-LSTM)	95.71%	94.52%	94.05%	94.29%	97.38%	95.59%
Test Accuracy (Inception-GRU)	96.00%	96.00%	96.67%	96.33%	97.33%	96.47%
Test Accuracy (Xception-LSTM)	88.00%	89.78%	90.45%	88.54%	89.99%	89.55%
Test Accuracy (Xception-GRU)	91.98%	91.02%	90.00%	91.90%	91.90%	91.36%

TABLE II

ACHIEVED ACCURACIES OF ALL MODELS ACROSS ALL FOLDS.

As illustrated in Table 2, the Inception-LSTM model outperforms other models, including Inception-GRU and Xception-LSTM. The Inception-LSTM’s average test accuracy of

95.59% indicates its superior ability to capture multi-scale features, with Xception-GRU showing the lowest average test accuracy at 91.36%. The deeper architecture and inception modules of Inception-v3 enhance spatial feature capture, while the LSTM's temporal capabilities contribute to higher overall performance, making Inception-LSTM highly effective for dynamic gesture recognition.

#### E. Future Work

Future work on dynamic hand gesture recognition could address several key areas to further enhance model performance. First, implementing advanced data augmentation techniques, such as geometric transformations and noise addition, would allow the model to generalize better across varied backgrounds and lighting conditions. This will improve robustness, especially in real-world applications where environments are less controlled. Second, transfer learning with domain-specific datasets, such as those specific to human-computer interaction, could enhance model adaptability and reduce training time when faced with smaller or unique gesture datasets.

Another promising direction is to incorporate lightweight CNN architectures and optimized recurrent units, such as GRUs, to make the model more suitable for mobile and embedded applications without compromising accuracy. This would expand the model's applicability to devices with limited computational resources, such as smartphones or wearable technology, for real-time gesture-based control. Finally, experimenting with attention mechanisms or transformer-based architectures could further improve the recognition accuracy by focusing on important temporal and spatial features within gesture sequences. This approach could especially help in distinguishing between visually similar gestures, such as "zoom-in" and "zoom-out." Overall, these enhancements aim to make gesture recognition systems not only more accurate but also more versatile and deployable in a wider range of interactive applications.

### V. CONCLUSION

In this work, dynamic hand gestures recognition is revisited using deep learning techniques for accuracy and precision. We also implement several models such as Inception-GRU, Xception-LSTM, and Xception-GRU that can capture both spatial and temporal features from gesture sequences. These CNN-RNN hybrids turned out to be effective in detecting subtle hand movements, recording a mean accuracy of 95.59% in real-time. Early stopping was employed to avoid overfitting, while robustness of the model was insured through 5-fold cross-validation. Our results considerably improve upon many related works, especially in the more realistic settings when dealing with dynamic gestures. Other works often deal only with static gestures or smaller datasets, yet our approach could perform comparably even under difficult conditions. Future work will focus on fine-tuning these models for improved performance, optimizing the architecture for real-time operation on mobile devices with limited computational power, and expanding the dataset to enhance robustness across diverse

environments. It may find applications in human-computer interaction, sign language interpretation, and gesture-based control systems.

### REFERENCES

- [1] S. Shin and W.-Y. Kim, "Skeleton-Based Dynamic Hand Gesture Recognition Using a Part-Based GRU-RNN for Gesture-Based Interface," *IEEE Access*, vol. 8, pp. 50236–50243, 2020.
- [2] A. S. M. Miah, M. A. M. Hasan, and J. Shin, "Dynamic Hand Gesture Recognition Using Multi-Branch Attention-Based Graph and General Deep Learning Model," *IEEE Access*, vol. 11, pp. 4703–4716, 2023.
- [3] W. Zhang and J. Wang, "Dynamic hand gesture recognition based on 3D convolutional neural network models," in *Proc. 2019 IEEE 16th Int. Conf. Networking, Sensing and Control (ICNSC)*, Banff, AB, Canada, 2019, pp. 363–368. doi: 10.1109/ICNSC.2019.8743159.
- [4] C. Linqin, C. Shuangjie, X. Min, Y. Jimin, and Z. Jianrong, "Dynamic hand gesture recognition using RGB-D data for natural human-computer interaction," *Journal of Intelligent & Fuzzy Systems*, vol. 32, no. 5, pp. 3495–3507, April 2017. DOI:10.3233/JIFS-169287.
- [5] M. U. Rehman, F. Ahmed, M. A. Khan, U. Tariq, F. Alfouzan, N. Alzahrani, and J. Ahmad, "Dynamic Hand Gesture Recognition Using 3D-CNN and LSTM Networks," *Computers, Materials & Continua*, September 2021. DOI:10.32604/cmc.2021.019586.
- [6] B. K. Triwijoyo, L. Y. R. Karnaen, and A. Adil, "Deep Learning Approach for Sign Language Recognition," *Jurnal Ilmiah Teknik Elektro Komputer dan Informatika*, vol. 9, no. 1, pp. 12–21, March 2023. DOI:10.26555/jiteki.v9i1.25051.
- [7] Z. Ren, J. Meng, and J. Yuan, "Depth camera based hand gesture recognition and its applications in Human-Computer-Interaction," presented at the *International Conference on Information, Communications and Signal Processing (ICICS)*, December 2011. DOI:10.1109/ICICS.2011.6173545.
- [8] D. Lu, C. Qiu, and Y. Xiao, "Temporal Convolutional Neural Network for Gesture Recognition," presented at the *2018 IEEE/ACIS 17th International Conference on Computer and Information Science (ICIS)*, Singapore, 6–8 June 2018. DOI:10.1109/ICIS.2018.8466467.
- [9] K. Bousbai and M. Merah, "A Comparative Study of Hand Gestures Recognition Based on MobileNetV2 and ConvNet Models," presented at the *2019 6th International Conference on Image and Signal Processing and their Applications (ISPA)*, Mostaganem, Algeria, 24–25 November 2019. DOI:10.1109/ISPA48434.2019.8966918.
- [10] A. S. M. Miah, M. A. M. Hasan, and J. Shin, "Dynamic Hand Gesture Recognition Using Multi-Branch Attention Based Graph and General Deep Learning Model," *IEEE Access*, vol. 11, pp. 4703–4716, January 9, 2023. DOI:10.1109/ACCESS.2023.3235368.
- [11] X. Tang, Z. Yan, J. Peng, B. Hao, H. Wang, and J. Li, "Selective spatiotemporal features learning for dynamic gesture recognition," *Expert Systems with Applications*, vol. 169, 1 May 2021, Article 114499. DOI:10.1016/j.eswa.2020.114499.
- [12] S. Jeeru, A. K. Sivapuram, D. G. León, J. Gröli, S. R. Yeduri, and L. R. Cenkaramaddi, "Depth Camera Based Dataset of Hand Gestures," *Data Brief*, vol. 45, Dec. 2022, Art. no. 108659. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2352340922008642>.
- [13] D. R. T. Hax, P. Penava, S. Krodell, L. Razova, and R. Buettner, "A Novel Hybrid Deep Learning Architecture for Dynamic Hand Gesture Recognition," *IEEE Access*, vol. 12, pp. 28761–28773, Feb. 2024. DOI:10.1109/ACCESS.2024.3365274.
- [14] D. Wang, C. Li, M. Chen, and L. Zhang, "Gesture recognition in complex environments using CNNs and LSTMs," *Pattern Recognition*, vol. 102, pp. 107241, 2023. DOI:10.1016/j.patcog.2022.107241.
- [15] S. Harris and A. Kumar, "Robust spatiotemporal modeling for dynamic gesture recognition," *IEEE Transactions on Multimedia*, vol. 25, pp. 3601–3613, 2024. DOI:10.1109/TMM.2024.3295471.