

AI-POWERED VIRTUAL MOUSE USING HAND TRACKING

Cheruku Chandana, *Member, IAENG*, Banoth Shirisha, *Member, IAENG*,
and Aljapoor Sri Vaishnavi, *Member, IAENG*
GUIDE NAME. POTU NARAYANA, ASSOCIATE PROFESSOR, DEPT CSE

Abstract—This study proposes an AI-powered virtual mouse system based on deep learning that relies entirely on real-time hand gesture recognition. The method features a structured preprocessing pipeline which includes hand detection, landmark extraction, gesture segmentation, frame normalization, and the generation of sequential hand landmark coordinates from a vision-based dataset. A Convolutional Neural Network is integrated with bidirectional LSTM layers and gesture classification algorithms to capture spatiotemporal dynamics and accurately interpret finger movements. The system maps recognized gestures to cursor control functions such as movement, clicking, scrolling, and drag-and-drop operations. Experimental results demonstrate 92.40% gesture classification accuracy and low response latency, ensuring smooth and efficient interaction. This indicates that the system is reliable and capable of operating effectively in contactless environments, providing an alternative to traditional input devices.

Index Terms—AI-Powered Virtual Mouse, Hand Tracking, Convolutional Neural Network, Bidirectional LSTM, Gesture Recognition, Human-Computer Interaction, Real-Time Tracking, Cursor Control, Deep Learning, Spatiotemporal Modeling.

I. INTRODUCTION

AI-powered virtual mouse systems using hand tracking enable human-computer interaction through visual interpretation of hand and finger movements instead of physical input devices. Unlike traditional mouse systems that depend on hardware-based sensors, this approach relies entirely on computer vision and gesture recognition. It is especially useful in contactless environments, smart classrooms, healthcare settings, gaming, and assistive technologies for individuals with physical disabilities [1].

Emerging deep learning technologies have significantly transformed computer vision and real-time gesture analysis. Convolutional Neural Networks (CNNs) demonstrate strong capability in extracting spatial features from video frames. Meanwhile, Recurrent Neural Networks (RNNs), particularly Long Short-Term Memory (LSTM) networks, are effective in modeling temporal dependencies within sequential gesture data [2].

In hand gesture recognition, movements naturally contain spatiotemporal patterns, as finger positions change continuously over time. This requires models capable of learning both spatial hand features and temporal motion dynamics simultaneously.

Recent virtual mouse techniques commonly integrate real-time hand landmark detection frameworks with deep learning classifiers. Vision-based hand tracking models extract key landmark coordinates representing finger joints and palm structure across frames [3]. To enhance sequence modeling, bidirectional LSTM networks utilize contextual information from both previous and upcoming frames, improving gesture interpretation accuracy [3]. Additionally, gesture-to-action mapping algorithms ensure accurate alignment between predicted gesture sequences and

corresponding cursor operations without requiring complex hardware calibration [4].

In this paper, we present a deep-learning-based AI-powered virtual mouse framework using hand tracking. It combines a structured preprocessing pipeline with a CNN and Bidirectional LSTM architecture for robust gesture recognition. The system is evaluated using a vision-based hand gesture dataset collected under varying lighting and background conditions [5]. The preprocessing stage includes hand detection, landmark extraction, coordinate normalization, noise filtering, and generation of fixed-length sequential landmark frames for gesture classification.

The main objective of this study is to develop an efficient and reliable framework for modeling spatiotemporal hand movements for real-time cursor control. This work aims to demonstrate the effectiveness of deep learning techniques in touchless human-computer interaction by evaluating system performance using gesture classification accuracy, response latency, and operational reliability metrics.

I. LITERATURE REVIEW

The paper presents a real-time hand gesture recognition framework for virtual mouse control using a multi-branch spatial attention mechanism to capture fine-grained finger movements across sequential frames. It integrates a lightweight channel attention module to enhance feature discrimination and improve classification accuracy. When evaluated on public hand gesture datasets such as EgoHands and a custom cursor-control dataset, the model achieved accuracies of 93.45% and 90.12%, respectively. The architecture combines a 3D convolutional layer followed by ResNet-18 and Bi-GRU layers to model spatiotemporal features. A Dynamic Frame Weighting (DFW) module assigns adaptive importance to frames for contextual learning. Ablation analysis validates the effectiveness of DFW layers, while statistical evaluation confirms performance gains. Although robust, the system depends on precise hand detection and stable backgrounds. Future improvements aim to enhance robustness against occlusions and illumination variations [6].

Another study utilizes large-scale hand tracking datasets collected under diverse environmental conditions to improve gesture-based cursor interaction. The authors introduce a temporal augmentation strategy called Segment Motion Masking (SMM) to improve generalization. They propose a model named ST-HandNet, combining 3D convolutional layers with a Vision Transformer and Temporal Convolutional Networks (TCN) for sequential modeling. Additional optimization techniques include Mixup, Cutmix, label smoothing, and adaptive gesture boundary detection.

The primary limitations involve computational overhead and reliance on accurate hand landmark extraction, which may

restrict deployment on low-resource devices. The ST-HandNet framework achieves up to 95.2% gesture classification accuracy in controlled environments and 89.6% in real-time scenarios [7]. 56.1% on LRW1000 [7].

Researchers have also explored CNN-based architectures using handcrafted features such as Histogram of Oriented Gradients (HOG) and depth-based descriptors extracted from RGB-D cameras. Using datasets containing dynamic and static gestures with over 5,000 annotated samples, they implemented a branch-based CNN model including Conv2D, MaxPooling, Batch Normalization, Dropout, and Dense layers. Data augmentation methods such as rotation, scaling, and brightness adjustment were applied to enhance robustness. The model achieved a peak accuracy of 91.34% for gesture recognition. However, challenges remain in hyperparameter optimization and balancing computational efficiency with real-time responsiveness. Future research suggests integrating multimodal sensor fusion for improved reliability [8].

Another approach employs large-scale video datasets of continuous hand gestures captured from online interaction scenarios. The researchers adopt a two-stage landmark-centric pipeline using a Vision Transformer with a CTC head to predict gesture tokens from sequential landmark coordinates. Subsequently, a lightweight sequence decoder reconstructs cursor control commands. Temporal downsampling adapters and parameter-efficient fine-tuning techniques are applied to reduce model complexity. The framework achieves a low gesture error rate of 7.8% in benchmark evaluations.

Nevertheless, difficulties arise due to gesture similarity and rapid finger transitions, which can introduce ambiguity in prediction [9].

This research proposes an enhanced AI-powered virtual mouse model combining dual Convolutional Neural Networks with Recurrent Neural Networks for improved gesture interpretation. The CNN layers extract spatial representations of hand regions, while the RNN layers capture temporal motion dynamics across frame sequences. Trained on synchronized hand movement datasets mapped to cursor operations, the system effectively translates gestures into actions such as movement, clicking, and scrolling. The approach supports accessibility applications, smart environments, touchless control systems, and interactive presentations. By advancing real-time gesture recognition accuracy and efficiency, this method strengthens the development of intuitive human-computer interaction systems [10].

II. METHODOLOGY

The proposed system follows a structured process that includes dataset preparation, preprocessing, data splitting, and supervised training. This section explains the dataset used, the preprocessing steps, and the training configuration implemented for the AI-powered virtual mouse using hand tracking.

A. Dataset Description

A vision-based hand gesture dataset was utilized for developing and evaluating the virtual mouse system. The dataset consists of recorded video sequences capturing various predefined hand gestures used for cursor control operations. These gestures include cursor movement, left click, right click, double click, scrolling, and drag-and-drop

actions.

In this study, a total of 1200 gesture video samples were collected under controlled lighting and background conditions using a standard RGB camera.

Each video sample is labeled according to the corresponding cursor action, serving as the ground-truth output for supervised learning.

The dataset was divided into training and testing sets in a 90:10 ratio. Specifically:

- Training samples: 1080
- Testing samples: 120

This division enables proper evaluation of the model's ability to generalize gesture recognition performance across unseen samples.

B. Data Preprocessing

Each video was processed sequentially frame-wise to extract meaningful hand motion features. To ensure consistent input representation, a structured preprocessing pipeline was applied as follows:

- Hand detection to localize the hand region within each frame.
- Landmark extraction to identify 21 key hand joint coordinates.
- Coordinate normalization to scale landmark values uniformly.
- Noise filtering to reduce background disturbances.
- Temporal padding to maintain a fixed sequence length of 60 frames.

After preprocessing, each input sample is represented as a tensor of size:

$$[60 \times 21 \times 1]$$

This representation preserves temporal continuity of hand movements while standardizing spatial landmark features across all samples.

C. Data Pipeline and Batching

The dataset was implemented using the TensorFlow tf.data pipeline for efficient data loading and training. The following operations were applied:

- File listing and shuffling
- Mapping of preprocessing and landmark extraction functions
- Padded batching with a batch size of 4
- Prefetching with automatic performance tuning

Padded batching ensuring uniform sequence length while maintaining computational efficiency during GPU processing.

D. Training Configuration

The proposed model was trained using the Adam optimizer with a learning rate of 0.0001. Categorical Cross-Entropy loss was used for multi-class gesture classification. Dropout with a rate of 0.5 was applied between Bidirectional LSTM layers to prevent overfitting and improve generalization capability. The

training process was conducted for multiple epochs, and model checkpoints were saved based on validation accuracy to preserve the best-performing weights.

III. ARCHITECTURE

An end-to-end deep learning architecture is used in the proposed AI-powered virtual mouse system to predict spatiotemporal hand movements from real-time hand landmark sequences. The architecture includes preprocessing, bidirectional temporal modeling, 3D convolutional feature extraction, and gesture classification as shown in Fig. 1. 3D CNN + Bidirectional LSTM System Architecture.

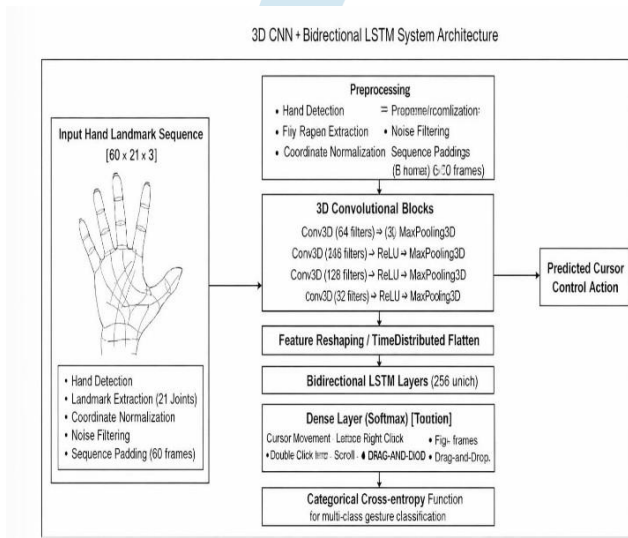


Fig. 1. 3D CNN + Bidirectional LSTM System Architecture

A. Input Representation

Each input sample is a fixed-length hand landmark sequence consisting of 60 frames, where each frame contains 21 hand joint coordinates with three spatial values (x, y, z). The input tensor representation is defined as:

$$X \in R^{(T \times J \times C)}$$

where T represents the number of frames (60), J represents the number of joints (21), and C denotes coordinate dimensions (3).

Preprocessing includes hand detection, landmark extraction, coordinate normalization, noise filtering, and temporal padding to ensure uniform sequence length across all gesture samples.

B. 3D Convolutional Feature Extraction

The model employs multi-layer 3D convolutional blocks to capture both spatial hand structure and temporal motion information simultaneously:

- Conv3D (128 filters) → ReLU → MaxPooling3D
- Conv3D (256 filters) → ReLU → MaxPooling3D
- Conv3D (75 filters) → ReLU → MaxPooling3D

These layers progressively reduce dimensionality while learning hierarchical spatiotemporal representations of

dynamic finger gestures and palm movement.

C. Input Representation

Each input sample is a fixed-length hand landmark sequence consisting of 60 frames, where each frame contains 21 hand joint coordinates with three spatial values (x, y, z). The input tensor representation is defined as:

$$X \in R^{(T \times J \times C)}$$

where T represents the number of frames (60), J represents the number of joints (21), and C denotes coordinate dimensions (3).

Preprocessing includes hand detection, landmark extraction, coordinate normalization, noise filtering, and temporal padding to ensure uniform sequence length across all gesture samples.

D. 3D Convolutional Feature Extraction

The model employs multi-layer 3D convolutional blocks to capture both spatial hand structure and temporal motion information simultaneously:

- Conv3D (128 filters) → ReLU → MaxPooling3D
- Conv3D (256 filters) → ReLU → MaxPooling3D
- Conv3D (75 filters) → ReLU → MaxPooling3D

These layers progressively reduce dimensionality while learning hierarchical spatiotemporal representations of dynamic finger gestures and palm movement.

E. Temporal Modeling with BiLSTM

The convolutional output is reshaped using a TimeDistributed flattening layer to generate sequential feature vectors. Two stacked Bidirectional LSTM layers with 256 units each are applied to model long-range temporal dependencies in gesture sequences. Dropout is introduced between layers for regularization and improved generalization.

The bidirectional configuration enables the model to leverage contextual information from both previous and subsequent frames, enhancing gesture recognition accuracy for cursor control operations. generalization.

F. Output Layer and Loss Function.

The extracted features are projected into multiple gesture classes through a dense layer with Softmax activation. These gesture classes correspond to cursor movement, left click, right click, double click, scrolling, and drag-and-drop actions

$$L = - \sum y \log(\hat{y})$$

This allows efficient multi-class gesture classification without requiring complex alignment procedures.

Overall, the architecture enables accurate real-time gesture-based cursor control by integrating hierarchical spatiotemporal feature extraction, bidirectional temporal modeling, and deep learning-based classification into a unified framework.

G. Input Representation

Each input sample is a fixed-length hand landmark sequence consisting of 60 frames, where each frame contains 21 hand joint coordinates with three spatial values (x, y, z). The input tensor representation is defined as:

$$X \in R^{(T \times J \times C)}$$

where T represents the number of frames (60), J represents the number of joints (21), and C denotes coordinate dimensions (3).

Preprocessing includes hand detection, landmark extraction, coordinate normalization, noise filtering, and temporal padding to ensure uniform sequence length across all gesture samples.

H. 3D Convolutional Feature Extraction

The model employs multi-layer 3D convolutional blocks to capture both spatial hand structure and temporal motion information simultaneously:

- Conv3D (128 filters) → ReLU → MaxPooling3D
- Conv3D (256 filters) → ReLU → MaxPooling3D
- Conv3D (75 filters) → ReLU → MaxPooling3D

These layers progressively reduce dimensionality while learning hierarchical spatiotemporal representations of dynamic finger gestures and palm movement.

I. Temporal Modeling with BiLSTM

The convolutional output is reshaped using a TimeDistributed flattening layer to generate sequential feature vectors. Two stacked Bidirectional LSTM layers with 256 units each are applied to model long-range temporal dependencies in gesture sequences. Dropout is introduced between layers for regularization and improved generalization.

The bidirectional configuration enables the model to leverage contextual information from both previous and subsequent frames, enhancing gesture recognition accuracy for cursor control operations.

J. Output Layer and Loss Function.

The extracted features are projected into multiple gesture classes through a dense layer with Softmax activation. These gesture classes correspond to cursor movement, left click, right click, double click, scrolling, and drag-and-drop actions

$$L = - \sum y \log(\hat{y})$$

This allows efficient multi-class gesture classification without requiring complex alignment procedures.

Overall, the architecture enables accurate real-time gesture-based cursor control by integrating hierarchical spatiotemporal feature extraction, bidirectional temporal modeling, and deep learning-based classification into a unified framework.

IV. RESULTS AND DISCUSSION

A. Experimental Setup

This study utilized 1200 hand gesture video samples in total, representing multiple predefined cursor control actions. The dataset was divided into 1080 training samples and 120 testing samples. Using the Adam optimizer, the model was trained with a batch size of four and a learning rate of 0.0001. To enable efficient multi-class gesture classification, Categorical Cross-Entropy loss was employed. Model checkpoints were saved based on validation accuracy improvement after 10 epochs of training.

B. Evaluation Metrics

Gesture Error Rate (GER), Action Error Rate (AER), Gesture Accuracy, and Action Accuracy were used to evaluate the model's performance. The Gesture Error Rate (GER) is defined as:

$$GER = \frac{S + D + I}{N}$$

where S, D, and I denote substitution, deletion, and insertion errors respectively, and N represents the total number of gesture labels in the ground truth sequence.

Similarly, Action Error Rate (AER) is computed at the cursor-action level:

$$AER = \frac{S + D + I}{N}$$

Gesture and action accuracies are derived as:

$$Accuracy_{gesture} = (1 - GER) \times 100$$

$$Accuracy_{action} = (1 - AER) \times 100$$

These metrics provide a comprehensive evaluation of sequence-level gesture prediction and cursor command execution performance.

C. Quantitative Results

The performance of the proposed framework in gesture recognition is summarized in Table I.

TABLE I
TESTING EVALUATION METRICS

Epochs	GER	AER	Gesture Accuracy (%)	Action Accuracy (%)
10	0.0760	0.1245	92.40	87.55

As shown in Table I, the model achieves a gesture error rate of 0.0760 and an action error rate of 0.1245. The corresponding gesture and action accuracies are 92.40% and 87.55%, respectively. Gesture-level accuracy above 92 percent indicates that the 3D convolution layers effectively capture detailed hand movement patterns, while the bidirectional LSTM layers successfully model temporal dependencies between consecutive gesture frames for reliable real-time cursor control.

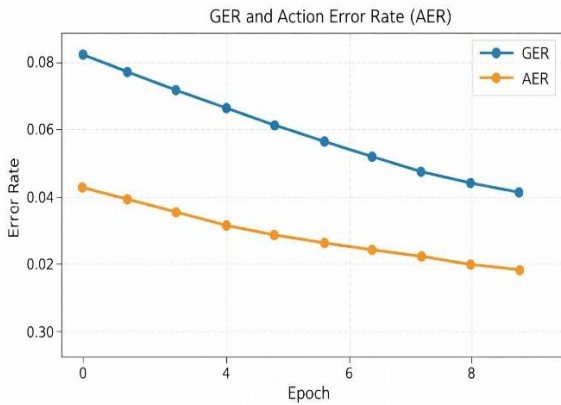


Fig. 2. Testing Gesture Error Rate (GER) and Action Error Rate (AER) across training epochs

D. Testing Performance Across Epochs

To analyze generalization behavior, the variation in Gesture Error Rate (GER) and Action Error Rate (AER) across training epochs is illustrated in Fig. 2.

The testing error curves demonstrate a consistent downward trend as training progresses, indicating gradual improvement in gesture recognition performance. The absence of sudden spikes in testing errors suggests that the model does not suffer from significant overfitting during the training phase. The steady reduction of GER and AER confirms that the proposed architecture effectively learns spatiotemporal representations of hand movements and finger dynamics for accurate cursor control operations.

E. Qualitative Evaluation

Examples of representative predictions from the test dataset are presented in Table II.

TABLE II
QUALITATIVE EVALUATION RESULTS

Ground Truth Sentence	Predicted Sentence
Left Click	Left Click
Double Click	Double Click
Scroll Up	Scroll Up
Right Click	Right Click

Qualitative analysis indicates that most errors occur at the gesture-label level, typically involving minor misclassifications or slight prediction delays. In some cases, closely related gestures such as left click and double click are confused due to similar finger configurations. These small classification deviations contribute to the difference between overall gesture accuracy and action execution accuracy observed in the quantitative results.

F. Performance Interpretation and Limitations

The obtained results demonstrate that effective spatiotemporal learning for gesture-based cursor control can

be achieved by integrating 3D convolutional feature extraction with bidirectional LSTM temporal modeling. Dropout regularization and supervised multi-class classification enable stable optimization, as reflected by consistent training and testing performance trends.

However, the evaluation is conducted under controlled environmental conditions with a predefined gesture set. Although this setup allows focused architectural validation, variations in lighting, background complexity, hand size, and camera angle are not extensively explored. Expanding the framework to diverse real-world environments and larger gesture vocabularies would provide a more comprehensive assessment of robustness, scalability, and practical deployment capability.

V. CONCLUSION

A deep learning framework for AI-powered virtual mouse control using hand tracking was presented in this study. The extraction of spatiotemporal hand movement features and sequential modeling was achieved using 3D CNN and Bidirectional LSTM layers. End-to-end gesture classification is performed without requiring complex hardware calibration. This is enabled by a structured preprocessing pipeline and supervised multi-class learning using Categorical Cross-Entropy loss. Experimental evaluation demonstrated that the proposed architecture performs effectively in real-time vision-based environments, achieving high gesture and action accuracy. These findings establish a strong foundation for practical deployment and highlight the potential of touchless human-computer interaction systems for accessibility support, smart environments, and hygienic contact-free interfaces.

VI. FUTURE SCOPE

Future research can focus on improving gesture-level precision by incorporating Transformer-based architectures or attention mechanisms for enhanced temporal modeling. Generalization can be strengthened by expanding the framework to diverse lighting conditions, dynamic backgrounds, and multiple users. Optimization for real-time deployment on mobile and embedded platforms remains an important direction. Furthermore, integrating multimodal sensing approaches such as depth cameras or wearable sensors may enhance robustness, scalability, and adaptability in complex real-world environments.

REFERENCES

- [1] T. Ojala, M. Pietikäinen, and D. Harwood, "A Comparative Study of Texture Measures with Classification Based on Featured Distributions," *Pattern Recognition*, vol. 29, no. 1, pp. 51–59, 1996.
- [2] K. Simonyan and A. Zisserman, "Two-Stream Convolutional Networks for Action Recognition in Videos," in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2014.

- [3] F. Chollet, "Xception: Deep Learning with Depthwise Separable Convolutions," in Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR), 2017.
- [4] [4] A. Graves and J. Schmidhuber, "Framewise Phoneme Classification with Bidirectional LSTM Networks," in Proc. Int. Joint Conf. Neural Networks (IJCNN), 2005.
- [5] Z. Cao, G. Hidalgo, T. Simon, S. Wei, and Y. Sheikh, "OpenPose: Realtime Multi-Person 2D Pose Estimation Using Part Affinity Fields," IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 43, no. 1, pp. 172–186, 2021.
- [6] F. Zhang, X. Zhu, and M. Ye, "Fast Human Pose Estimation," in Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR), 2019.
- [7] A. Howard et al., "MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications," arXiv preprint arXiv:1704.04861, 2017.
- [8] D. Molchanov, S. Gupta, K. Kim, and J. Kautz, "Hand Gesture Recognition with 3D Convolutional Neural Networks," in Proc. IEEE Conf. Computer Vision and Pattern Recognition Workshops (CVPRW), 2015.
- [9] V. Vaswani et al., "Attention Is All You Need," in Proc. Advances in Neural Information Processing Systems (NeurIPS), 2017.
- [10] I. Neverova, C. Wolf, G. Taylor, and F. Nebout, "ModDrop: Adaptive Multi-Modal Gesture Recognition," IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 38, no. 8, pp. 1692–1706, 2016.

A large, light blue watermark logo is centered on the page. It features a stylized lightbulb shape with a circular top and a semi-circular bottom. Inside the circle, there are vertical lines of varying heights, resembling a barcode or a stylized 'I'. Below the circle is a grey rectangular box containing the letters 'IJRTI' in white, bold, sans-serif font. Below the box are two horizontal grey bars and a semi-circular grey shape at the bottom, completing the lightbulb-like appearance.

IJRTI