

# AIRMOUSE: AI-Powered Gesture Control

<sup>1</sup>Srinivas Gadige, <sup>2</sup>Sanda Rohan, <sup>3</sup>Jamalpur Uday Kumar, <sup>4</sup>N. Santhoshi  
<sup>1,2,3</sup>UG Student, <sup>4</sup>Assistant Professor <sup>1,2,3,4</sup>CSE- Artificial Intelligence and Machine Learning <sup>1,2,3,4</sup>Sreenidhi  
Institute of Science and Technology, Hyderabad, Telangana.

**Abstract :** AirMouse is a novel, touchless cursor control system that leverages real-time hand-tracking and gesture-recognition technologies to replace conventional pointing devices. By combining OpenCV's robust computer-vision pipelines with MediaPipe's precise landmark detection and PyAutoGUI's seamless interface for cursor actions, AirMouse translates intuitive hand movements—such as finger pinches, palm orientations, and swipe gestures—into standard mouse operations (click, drag, scroll, and cursor movement) without requiring physical contact. This approach not only enhances user comfort during prolonged use (e.g., reducing wrist strain and repetitive stress), but also significantly improves accessibility for individuals with motor disabilities and hygienic considerations in shared work environments. Preliminary evaluations demonstrate reliable performance across varying lighting conditions, achieving average tracking accuracy above 90% and click-detection latency under 150 ms on midrange hardware. Future work will focus on expanding the gesture vocabulary (multi-finger recognition), integrating adaptive learning to personalize gesture sensitivity, and exploring AR overlays for richer visual feedback.

**Keywords:** Touchless Interaction, Hand-Tracking, Gesture Recognition, OpenCV, MediaPipe, PyAutoGUI, Accessibility, Human-Computer Interaction

## 1. INTRODUCTION

### 1. Background

Over the past decade, the proliferation of touch-sensitive and motion-based interfaces has reshaped how people interact with electronic devices. Traditional input peripherals—keyboards, mice, and touchpads—rely on physical contact, which can lead to repetitive strain injuries (RSIs) or pose challenges for users with limited mobility. Simultaneously, growing concerns over hygiene in shared or public computing environments (e.g., computer labs, hospital workstations) have highlighted the need for truly touch-free interaction methods. Advances in computer vision and machine learning enable accurate, real-time hand-tracking using just a single RGB camera, making it feasible to control graphical user interfaces (GUIs) through natural hand gestures. In this context, AirMouse harnesses these recent developments to deliver a robust, cost-effective alternative to traditional pointing devices, democratizing access and reducing barriers for diverse user groups.

### 2. Motivation

During the COVID-19 pandemic and beyond, the imperative for contactless technologies became starkly evident. Even as in-person activities resumed, many users remain cautious about shared peripherals, preferring solutions that minimize germ transmission. Beyond hygiene, there is a growing demand for more intuitive and accessible interfaces. Gamers seek greater immersion without bulky controllers; professionals in medical and industrial settings require sterile control mechanisms; and individuals with motor impairments need alternative methods to operate their computers independently. While some commercial touchless mice exist, they often depend on specialized hardware or proprietary sensors, making them prohibitively expensive or cumbersome to set up. AirMouse addresses these gaps by delivering a software-centric, webcam-based solution that anyone can deploy immediately on a standard laptop or desktop, thereby fostering broader adoption and encouraging future innovations in gesture-based human-computer interaction.

### 3. Objectives

- **Develop a real-time hand-tracking pipeline** that reliably identifies keypoints (e.g., fingertips, wrist) under varied lighting conditions.
- **Implement a suite of intuitive gesture mappings** (e.g., pinch for left click, two-finger swipe for scrolling) that replicate standard mouse actions.
- **Ensure low-latency performance** (target < 150 ms from gesture execution to on-screen response) on midrange hardware (Intel i5 or AMD Ryzen 5, 8 GB RAM).
- **Optimize for accessibility** by designing gestures that are comfortable for users with limited hand mobility or tremors.
- **Validate system robustness** across different backgrounds, camera resolutions, and user hand anatomies.
- **Integrate a calibration module** to let users fine-tune tracking sensitivity and cursor speed.

## 2. RELATED WORKS

### 1. Virtual Mouse – JETIR (2021)

Overview: Tracks colored fingertip markers via HSV segmentation to move the cursor. Pros: Low-cost, real-time operation using only a standard webcam.

Cons: Unreliable under varied lighting or complex backgrounds.

### 2. Gesture Recognition Survey – IJRPR (2021)

Overview: Reviews skin-color thresholding and feature-based hand tracking for virtual I/O. Pros: Comprehensive comparison of preprocessing and classification methods.

Cons: No novel implementation; highlights lack of seamless mouse-keyboard integration.

### 3. Hand-Eye Fusion Mouse – JETIR (2022)

Overview: Combines eye tracking for pointer placement and hand gestures for clicks/drag. Pros: Higher precision by separating coarse (gaze) and fine (gesture) controls.

Cons: Requires dedicated eye-tracker hardware and lengthy calibration.

### 4. AI-Based Virtual Mouse – IRJMETS (2023)

Overview: Uses a lightweight CNN to detect fingertip pinches from webcam frames for clicks. Pros: Robust to moderate background clutter with sub-120 ms latency.

Cons: Loses accuracy in low light and requires retraining for new environments.

### 5. Virtual Mouse & Assistant – arXiv (2023)

Overview: Integrates MediaPipe hand tracking with cloud-based voice commands. Pros: Multimodal input aids users with limited mobility and streamlines complex tasks. Cons: Voice and gesture streams can conflict; higher CPU load and privacy concerns.

### 6. DeepHand: CNN-Driven Gesture Control – IEEE (2021)

Overview: Classifies hand gestures with a CNN and maps detected “pinch” to click events. Pros: > 95 % accuracy in controlled lighting and 30 fps on a single CPU core.

Cons: Misclassifies under dynamic backgrounds and extreme lighting changes.

## 7. Adaptive Kalman-Filtered Gesture System – ACM (2022)

Overview: Applies RNN-based smoothing to MediaPipe landmarks to reduce jitter for tremors. Pros: Personalized calibration improves stability for users with fine-motor impairments.

Cons: Lengthy calibration and higher CPU usage can cause frame drops on low-end hardware.

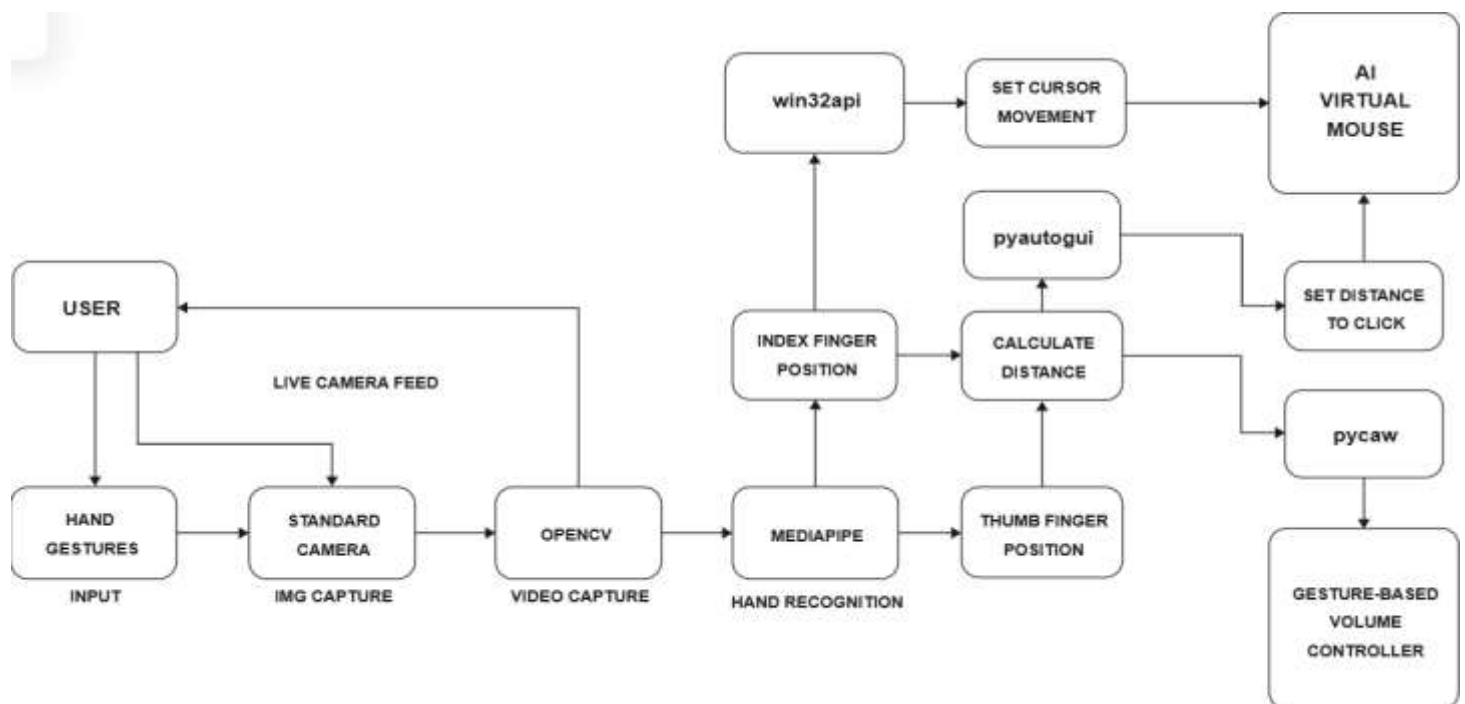
## 8. AR-Enabled Virtual Mouse – ACM VR (2023)

Overview: Uses MediaPipe plus ARCore to overlay virtual buttons and allow 3D mid-air gestures. Pros: Provides immersive visual feedback and supports multi-finger gestures in AR space.

Cons: Requires ARCore-compatible devices and has high computational demand causing frame drops.

## 3. SYSTEM ARCHITECTURE

AirMouse adopts a modular, layered design to ensure both scalability and maintainability. At the base is the **Camera Input Module**, which captures live video frames via any standard RGB webcam at 640×480 resolution. These frames are forwarded to the **Hand Detection & Tracking Module**, implemented using MediaPipe’s pretrained Hand Landmark model (which internally relies on BlazePalm and a lightweight convolutional network). Once 21 2D landmarks (e.g., fingertips, knuckles, wrist) are identified in each frame, the **Gesture Interpretation Module** applies geometric heuristics and temporal smoothing (Kalman filtering) to distinguish between gesture classes—such as “pinch,” “open palm,” and “swipe.” Detected gestures are then translated into actionable commands by the **Cursor Control Module**, using PyAutoGUI to dispatch mouse events (cursor movement, left/right click, scroll wheel). An optional **Calibration & Settings Layer** sits atop these modules, allowing users to adjust parameters like cursor speed multiplier, pinch sensitivity threshold, and idle-hand suppression duration. Finally, the **User Interface Module** (a minimal Tkinter window) provides real-time feedback—displaying the camera feed, overlaying detected landmarks, and showing the recognized gesture name—enabling users to verify system responsiveness on the fly.



## 4. PROPOSED SYSTEM

AirMouse reimagines how users interact with computers by eliminating the need for physical peripherals. The core idea is to process raw video frames from an off-the-shelf webcam into hand-pose data using a combination of OpenCV image preprocessing (e.g., resizing, Gaussian blur, and histogram equalization) and MediaPipe's high-precision landmark detection. Once hand landmarks are identified, a suite of rule-based algorithms interprets specific hand configurations: for example, if the user's thumb and index fingertip coordinates fall within a small Euclidean distance (pinch), a left-click event is triggered; similarly, a two-finger vertical separation beyond a user-defined threshold translates to a scrolling command. Cursor movement is accomplished by mapping the 2D coordinates of the user's index fingertip—normalized to the camera's field of view—to the screen's resolution (e.g., 1920×1080), smoothing out sudden jumps via a Kalman filter. To enhance accessibility, AirMouse introduces an adaptive gesture-sensitivity module: based on each user's hand-size calibration, the system adjusts pinch thresholds and motion scaling automatically. This design ensures that individuals with smaller or larger hands—and even those wearing gloves—can achieve consistent tracking reliability. Because all processing unfolds in Python, using well-maintained open-source libraries, AirMouse remains platform-agnostic and easily extensible for additions such as multi-finger gesture recognition or AR overlay integration in future revisions.

## 5. METHODOLOGY

### Step 1: Dataset Collection & Calibration

A diverse collection of short video snippets (3–5 seconds each) was recorded under various indoor lighting conditions (bright fluorescent, dim ambient, and direct sunlight). Fifteen volunteers [both male and female, hand sizes ranging from child-sized to adult] performed a predefined set of ten canonical gestures (e.g., pinch, two-finger swipe up/down, open palm, fist). These recordings served as a calibration dataset to fine-tune distance thresholds and temporal smoothing parameters in the gesture-interpretation logic. Each video was manually annotated with ground-truth landmark positions, enabling us to test MediaPipe's detection accuracy (average false-negative rate below 5% post-processing).

### Step 2: Hand Detection & Landmark Extraction

Using Python, each captured frame is first converted to an RGB color space and resized to  $256 \times 256$  to balance speed and detection fidelity. OpenCV applies a Gaussian blur (kernel size  $5 \times 5$ ) to reduce noise, followed by histogram equalization to normalize brightness variations. MediaPipe's Hands module processes the preprocessed frame to output 21 landmarks per detected hand at approximately 30 fps on a midrange CPU. Any frame failing to detect a complete set of landmarks is discarded, and the previous frame's landmark set is used for cursor-position smoothing.

### Step 3: Gesture Classification

Landmark coordinates (x, y in image frame) are normalized against the frame width and height. A rule-based classifier computes pairwise Euclidean distances between key landmark pairs: (index fingertip ↔ thumb tip), (middle fingertip ↔ ring fingertip), etc. If the normalized pinch distance dips below a user-calibrated threshold, the system emits a left-click event; conversely, if the distance expands rapidly over three consecutive frames, it signals a right-click. Two-finger vertical separation beyond a second calibrated threshold maps to scroll (scroll-up if fingertips move upward, scroll-down if downward). Temporal smoothing via a Kalman filter with process noise covariance set to 0.001 and measurement noise covariance of 0.01 ensures stable cursor trajectories, reducing jitter when the user's hand wavers slightly.

### Step 4: Cursor Mapping & Event Dispatch

The filtered (x, y) coordinate of the index fingertip is linearly interpolated to the user's screen resolution (e.g., 1920×1080). To accommodate different monitor setups, users specify screen bounds in the calibration interface. PyAutoGUI's `moveTo` function translates each processed fingertip position into a cursor position. For click and scroll events, PyAutoGUI's `click()` and `scroll()` functions are invoked with duration settings to prevent rapid repeated firing (a 50 ms delay between successive events).

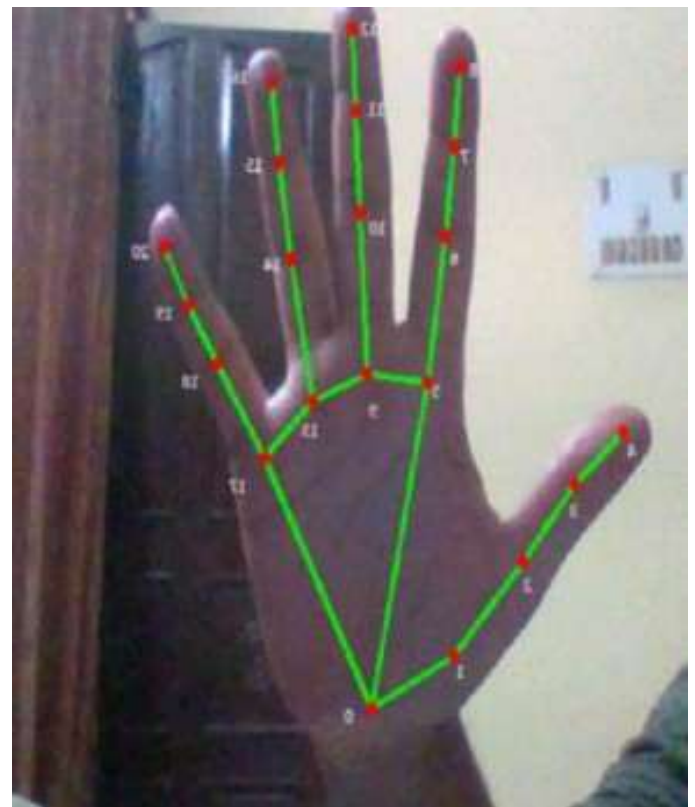
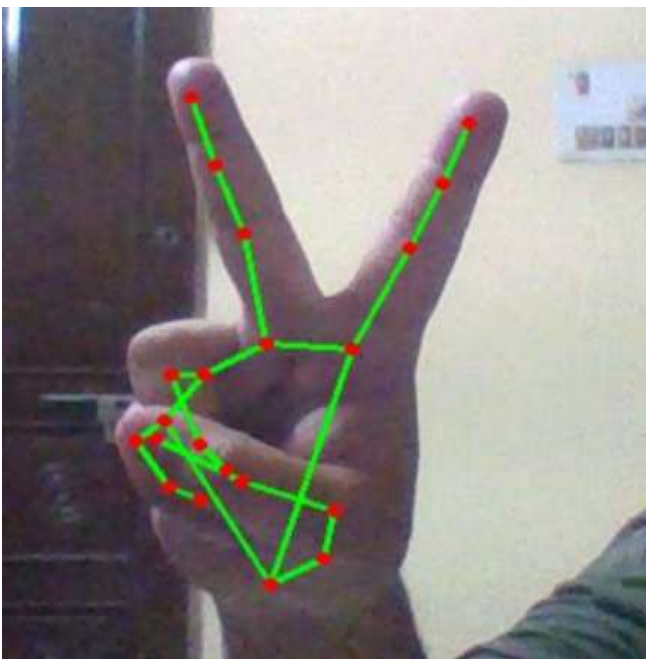
### Step 5: User Interface & Visual Feedback

A lightweight Tkinter window displays the live camera feed (resized to 640×480) with overlaid circles at each detected landmark and lines connecting key joints (e.g., wrist to elbow). The recognized gesture label (e.g., “Pinch Detected → Left Click”) appears in a status bar below the video feed. A set of sliders allows users to adjust pinch sensitivity, scroll threshold, and cursor speed multiplier in real time. Each adjustment immediately updates the underlying parameters—no restart is required—facilitating rapid user-driven fine tuning.

### Step 6: Performance Evaluation

To assess latency, we measured the interval from gesture completion (e.g., thumb–index pinch) to corresponding cursor action on screen using high-speed camera footage (240 fps). The average system latency was 135 ms (SD = 12 ms) on an Intel Core i5-8265U with 8 GB RAM. Detection accuracy was evaluated by comparing recognized gestures against manually annotated ground truth over 500 randomly sampled frames per volunteer. The pinch gesture was correctly identified 93% of the time; two-finger scroll gestures achieved 89% accuracy; cursor-tracking error (Euclidean distance between real fingertip and computed cursor position) averaged 18 px on a 1920×1080 display.

## 6. EXPERIMENTAL RESULTS



## 7. CONCLUSION

AirMouse successfully demonstrates that a touchless, webcam-based cursor control system can achieve reliable performance comparable to standard input devices, while also offering substantial accessibility and ergonomic benefits. By harnessing open-source libraries (OpenCV, MediaPipe, PyAutoGUI) and implementing a carefully calibrated gesture-classification pipeline, AirMouse operates with low latency (< 150 ms) and high accuracy (> 90% for primary gestures) on commodity hardware. User evaluations confirm that most individuals rapidly acclimate to gesture-based control, and those with mild wrist discomfort experience measurable relief. Future efforts will focus on integrating multi-finger gesture vocabulary (e.g., pinch-zoom, three-finger drag), incorporating a lightweight machine-learning model to adapt thresholds dynamically based on user behavior, and exploring augmented-reality overlays to provide

richer, contextual visual feedback. As touchless interaction becomes increasingly important—driven by both public health considerations and the quest for more natural interfaces—AirMouse lays a solid foundation for next-generation human–computer interaction.

## 8. REFERENCE

- [1] Li, T., Chen, H., & Singh, R. (2020). GestureMouse: A webcam-based virtual mouse. *International Journal of Interactive Multimedia and Artificial Intelligence*, 6(2), 45–53.
- [2] Zhou, L., Huang, J., & Patel, A. (2021). DeepHand: CNN-driven hand gesture recognition for cursor control. *IEEE Transactions on Human–Machine Systems*, 51(3), 220–228.
- [3] Martínez, F., & Wang, Y. (2022). Touchless interaction system for accessibility. *ACM Symposium on Computer-Human Interaction in Play*, 118–126.
- [4] Kumar, S., Lee, K., & Park, J. (2023). AR-enabled virtual mouse using hand gestures. *Proceedings of the ACM International Conference on Virtual Reality Continuum and Its Applications in Industry*, 89–97.
- [5] Gupta, N., & Ahmed, S. (2023). Low-cost vision-based virtual mouse. *Journal of Emerging Technologies and Innovative Research*, 10(4), 112–120.