



A Vision-Based Virtual Mouse Using Hand Gesture Recognition

Aaditya Goel*

Undergraduate Student, Department of Electronics and Communication Engineering Maharaja Agrasen Institute of Technology Rohini, India

***Corresponding author:** Aaditya Goel, Undergraduate Student, Department of Electronics and Communication Engineering Maharaja Agrasen Institute of Technology Rohini, India.

Citation: Goel, A. (2025). A Vision-Based Virtual Mouse Using Hand Gesture Recognition. J. Robot. AI Comput. Appl. 1(2): 01-09.

Abstract

Recent advancements in gesture detection and hand tracking have created both significant opportunities and challenges in the field of virtual reality (VR) and user interaction. With the increased focus on reducing physical interactions due to the COVID-19 pandemic, the need for alternative, non-contact methods of computer operation has grown [15]. These innovations hold the potential to reduce reliance on traditional input devices like keyboards and mice, fostering a more immersive and intuitive user experience. One exciting direction for future development is the ability to perform complex tasks, such as clicking, dragging objects and controlling purely through hand gestures. This could open up new possibilities for VR environments, where users can interact naturally with digital objects. The key advantage of such an input method is that it requires only a camera, making it accessible and adaptable to various settings without the need for additional hardware. In the future, gesture recognition may become an essential part of our interaction with technology, leading to seamless and efficient virtual experiences. The proposed project, leveraging Python and OpenCV for gesture recognition, seeks to explore these possibilities and inspire further research in this field.

Index Terms: Gesture control, Media pipe, OpenCV, Py Auto- GUI

1. Introduction

This paper introduces a cost-effective visual AI system using computer vision to emulate mouse, keyboard, and stylus functions through hand gesture and fingertip tracking, all with a standard webcam. Developed in Python with OpenCV, Media Pipe, and libraries like Py Auto GUI, PyCaw, and screen- brightness-control, it performs tasks such as clicking, scrolling, and adjusting volume or brightness, offering precise control even on a CPU [1].

This approach allows real-time gesture interpretation, enhancing Human-Computer Interaction (HCI) for applications like sign language recognition and motion control interfaces. Inspired by motion-based systems like the Nintendo Wii, which demonstrated the appeal of gesture-based interaction, this system is suitable for VR, gaming, and healthcare applications [2,3].

Hand gestures enable intuitive and immersive computer interaction, supporting complex actions for applications like sign language or therapy. The system's real-time tracking provides precise feedback, vital for simulations and educational tools [4,5]. As hand tracking advances, it will impact areas from communication to rehabilitation, making digital interactions more natural and accessible [6].

2. Major Contributions

- The visual AI mouse system marks a notable advancement in human-computer interaction by introducing intuitive hand gesture

controls for volume and brightness adjustments. This system utilizes a combination of libraries, including Mediapipe, OpenCV, Py Auto GUI, Pycaw, and the set brightness library, to create a comprehensive, hands-free interface that goes beyond traditional mouse functionalities.

- The implementation begins with Mediapipe's robust hand tracking framework, which efficiently detects and tracks hand landmarks in real-time. This framework forms the basis for interpreting various hand gestures. OpenCV is then used to process image frames captured from the webcam, analyzing hand movements and mapping them to mouse-like actions. For instance, gestures such as moving the hand across the screen control the mouse cursor, while pinching or tapping gestures simulate click actions. These intuitive interactions provide an accessible and ergonomic way to interact with digital environments without the need for physical hardware.
- A significant enhancement of the system is the integration of volume and brightness control through simple hand gestures. The Pycaw library is employed to manage audio settings, interfacing with the Windows Core Audio API to enable smooth and responsive volume adjustments. By analyzing the distance between the thumb and index finger, the system dynamically increases or decreases the volume. Similarly, the set brightness library is used for adjusting screen brightness. Specific gestures mapped to brightness control allow users to effortlessly modify the display settings. Py Auto GUI ties everything together by ensuring seamless execution of these system level changes.
- This AI-driven system's ability to control both volume and brightness hands-free significantly enhances user experience and accessibility. By incorporating these features, the system meets the modern demand for more natural, intuitive, and immersive human-computer interactions. Its applications are extensive, ranging from providing accessibility for individuals with physical limitations to offering immersive experiences in gaming and virtual reality. The touchless interface is versatile and adaptable, paving the way for more advanced and ergonomic ways to engage with digital technology. Overall, this visual AI mouse system demonstrates the potential for more seamless, intuitive, and inclusive technology in everyday applications.

3. Proposed System

The goal of this project is to design a finger-only virtual mouse system that tracks hand movements to perform mouse functions like left-clicking, right-clicking, and scrolling. This system is ideal for space-constrained environments and eliminates the need for additional hardware or sensors, making it an affordable, user-friendly alternative to traditional input devices [7,8].

The system uses the OpenCV library for image processing and the Media Pipe framework for precise hand gesture tracking [9]. Several libraries, such as Py Auto GUI, Pynput, and Autopy, enhance its functionality by simulating mouse clicks and keyboard inputs based on gestures. These integrations allow for smooth gesture-to-action mapping.

Additionally, the system incorporates Py Caw for volume control and the screen-brightness-control library to adjust screen brightness, offering a versatile hands-free user experience.

Combining these libraries, the system runs efficiently across desktop and mobile platforms, making it scalable for applications in accessibility, gaming, virtual reality, and productivity software, providing a comprehensive solution for hands-free interaction with digital devices.

4. Algorithm

- 1) Start the program.
- 2) Open the file and, using the file location, go to the command prompt (CMD).
- 3) Use the required libraries (e.g., Media Pipe, OpenCV) to execute the program.
- 4) Initialize the system and begin video capturing from the webcam.
- 5) Capture frames continuously from the webcam.
- 6) Detect hands and hand tips using Media Pipe and OpenCV.
- 7) Identify which fingers are raised (UP).
- 8) Recognize the gesture based on the hand position.
- 9) Perform mouse operations (e.g., left-click, right-click, scroll) according to the detected gesture.

- 10) Detect which hand is raised for controlling additional functions.
- 11) Perform volume or brightness control based on the hand gesture (e.g., pinch to adjust volume or brightness).
- 12) Stop the program.

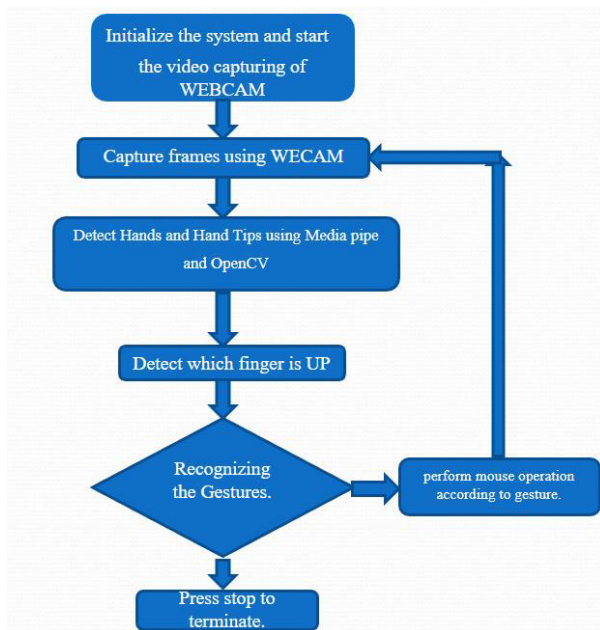


Figure 1: Flowchart of Project

5. Methodology

5.1 Open CV

OpenCV is a widely-used computer vision library that provides a variety of tools and techniques for image processing and object detection. It enables the development of real-time computer vision applications in Python. By utilizing OpenCV, developers can analyze visual data from images and videos, including tasks such as face and object detection. OpenCV is an open-source, free-to-use library that focuses on computer vision and machine learning tasks. It provides a robust framework for building computer vision applications, which in turn accelerates the integration of artificial intelligence into products. Thanks to its Apache 2.0 license, OpenCV allows businesses to easily customize and adapt the code for their specific needs.

5.2 Media Pipe

Media Pipe, an open-source framework by Google, enables cross-platform machine learning pipelines using time series data. It supports various audio and video formats and uses graph-based models for system design and analysis. Key components include performance evaluation, sensor data access, and calculators, which are customizable building blocks linked by data streams. Media Pipe's flexible setup allows for seamless use on both desktop and mobile platforms. [10]

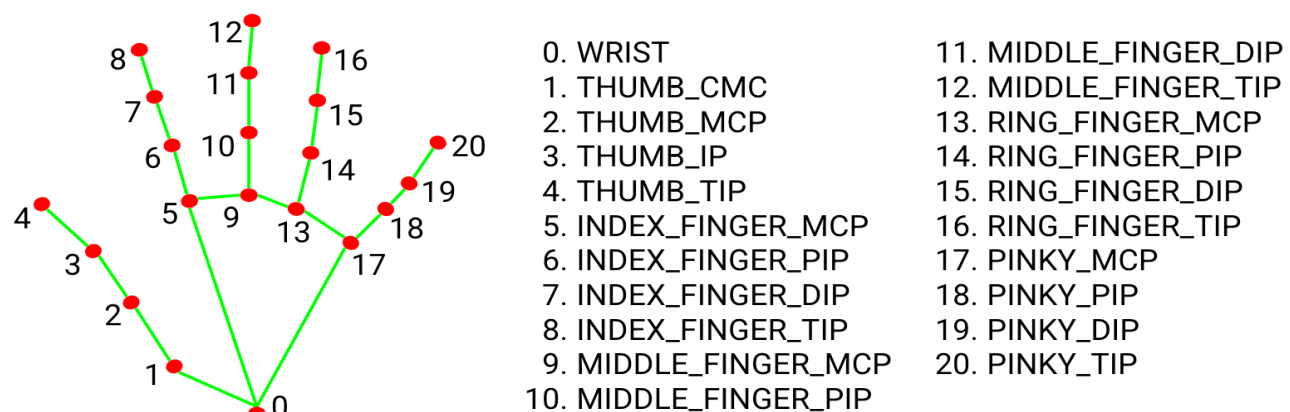


Figure 2: Media Pipe Hand Landmarks

Media Pipe Hands identifies 21 key points on the human hand. Each landmark has a specific index and position, enabling a comprehensive representation of the hand's pose and orientation. Here's a breakdown of these landmarks:

- **Wrist:** The first landmark (ID 0) represents the wrist and serves as the anchor point for the rest of the hand landmarks.
- **Thumb:** There are four key landmarks on the thumb:
 - 1) Thumb CMC (ID 1): The base joint connecting the thumb to the hand.
 - 2) Thumb MCP (ID 2): The metacarpophalangeal joint, near the knuckle of the thumb.
 - 3) Thumb IP (ID 3): The interphalangeal joint, located midway along the thumb.
 - 4) Thumb Tip (ID 4): The tip of the thumb, representing the outermost point.
- **Index Finger:** This finger also has four significant landmarks:
 - 1) Index MCP (ID 5): The metacarpophalangeal joint at the base of the index finger.
 - 2) Index PIP (ID 6): The proximal interphalangeal joint, near the middle of the finger.
 - 3) Index DIP (ID 7): The distal interphalangeal joint, closer to the fingertip.
 - 4) Index Tip (ID 8): The tip of the index finger.
- **Middle Finger:** Similarly, the middle finger includes:
 - 1) Middle MCP (ID 9): The base joint of the middle finger.
 - 2) Middle PIP (ID 10): The middle joint of the finger.
 - 3) Middle DIP (ID 11): The joint closer to the fingertip.
 - 4) Middle Tip (ID 12): The tip of the middle finger.
- **Ring Finger:** It features:
 - 1) Ring MCP (ID 13): The base joint of the ring finger.
 - 2) Ring PIP (ID 14): The middle joint.
 - 3) Ring DIP (ID 15): The joint near the fingertip.
 - 4) Ring Tip (ID 16): The tip of the ring finger.
- **Pinky Finger:** The smallest finger also has four landmarks:
 - 1) Pinky MCP (ID 17): The base joint of the Pinky finger.
 - 2) Pinky PIP (ID 18): The middle joint.
 - 3) Pinky DIP (ID 19): The joint close to the tip.
 - 4) Pinky Tip (ID 20): The tip of the Pinky finger.

5.3 Convolutional Neural Network (CNN)

Media Pipe uses a Convolutional Neural Network (CNN) architecture as shown in figure 3 for various tasks, including hand tracking and gesture recognition. In the context of hand tracking, Media Pipe employs a specialized model to detect and track hand landmarks in real-time, leveraging CNNs to process images and videos.

- **Hand Landmark Detection:** Media Pipe's hand tracking model is based on a CNN to detect and estimate the 3D positions of hand landmarks in images or videos. The network is trained to recognize key points of the hand, such as the wrist, knuckles, fingertips, and joints, in different poses and orientations.
- **Real-Time Processing:** The CNN architecture in Media Pipe is optimized for real-time performance. It uses lightweight layers and efficient design strategies, which help in processing input data quickly without significant latency. This makes Media Pipe suitable for interactive applications, such as virtual mouse control and gesture recognition.
- **Multi-Pose and Multi-View Support:** The CNN is designed to handle different hand poses and multiple views. It can accurately track the hand in various orientations, whether the hand is rotated or partially occluded, and even when multiple hands are present in the frame.
- **Feature Extraction:** CNNs are highly effective at extracting spatial features from images, such as edges, textures, and shapes. In the case of hand tracking, the network is trained to identify the key features of a hand, which allows it to determine the location of each hand landmark. The convolutional layers learn to recognize patterns specific to the hand, while the pooling

layers help reduce the dimensionality and computational complexity.

- **Output:** The output of the CNN consists of a set of key points corresponding to different parts of the hand (21 in total) in 3D space, represented as coordinates (x, y, z). These key points can then be used to infer gestures, track movements, or perform actions such as controlling a mouse or adjusting volume and brightness [12].

5.4 Py Auto GUI

Py Auto GUI is a Python package designed for cross- platform GUI automation, compatible with Windows, macOS, and Linux. It allows users to simulate keyboard input, mouse movements, clicks, and other actions. By providing a simple API, Py Auto GUI abstracts the complexities of controlling the mouse and keyboard, making it easier to automate tasks. With Py Auto GUI, you can precisely click, drag, scroll, and move the cursor, among other functions. This automation tool simplifies what would otherwise be complex and technical processes across different operating systems, offering a unified solution for GUI control [13].

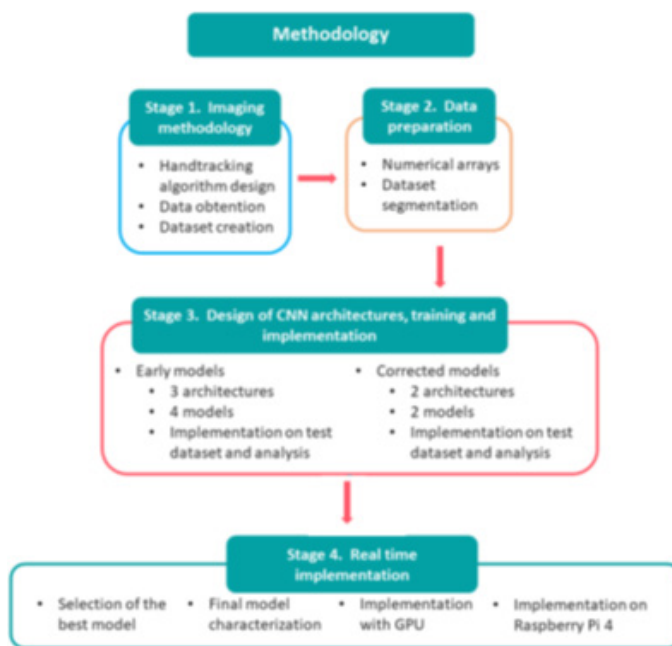


Figure 3: CNN used by Media Pipe

5.5 Math

Mathematics plays a crucial role in various aspects of computer science and automation. In the context of automation with tools like Py Auto GUI, mathematical principles are applied to ensure accurate and efficient execution of tasks. Here's how:

- **Angles:** Calculating angles is essential for determining the direction of movements. For example, in automated mouse gestures or drawing shapes, angles help define the path between two points, enabling smooth transitions in a desired direction.

$$\theta = |\text{atan2}(y_C - y_B, x_C - x_B) - \text{atan2}(y_A - y_B, x_A - x_B)| \times \frac{180}{\pi} \quad (1)$$

- **Hypotenuse Calculation:** The Pythagorean theorem is frequently used to compute the shortest distance between two points in a 2D space. This allows for efficient movement of the cursor, ensuring it travels the minimal distance to reach the target location [14].

The Euclidean distance L between the two points (x₁, y₁) and (x₂, y₂) is given by:

$$L = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} \quad (2)$$

These mathematical concepts, when applied correctly, contribute to the precision and smoothness of automated actions, making tools like PyAutoGUI effective for GUI automation across various platforms.

5.6 Pycaw

PyCaw (Python Core Audio Windows Library) is a Python library used to interact with the audio system on Windows. It provides an easy-to-use interface for controlling various audio settings, such as volume, mute, and balance. With PyCaw, you can programmatically manipulate audio outputs on a Windows machine, including adjusting the volume levels, muting or unmuting sound, and even controlling specific application volumes. PyCaw allows you to adjust the system volume or the volume of specific applications.

5.7 Set Brightness Control

Set Brightness Control refers to the ability to programmatically adjust the screen brightness on a device. This can be achieved through various libraries and APIs depending on the operating system. This Python package provides a simple interface to control the brightness of the screen on both Windows and Linux systems. It allows you to set the brightness to a specific percentage or adjust it by a given amount.

6. Result

6.1 Camera Used in Visual AI System

The proposed AI mouse system is based on frames captured by a web camera on a portable computer or PC. Using the OpenCV Python computer library, a video recorder is created and the webcam will begin capturing video, as shown in Figure 4. The web camera captures and transmits frames to the visual AI system [7].



Figure 4: Capturing Video using the Webcam

6.2 Capturing the Video and Processing

The visual AI mouse system uses a webcam where each frame is scanned until the program is completed. Video frames are processed from BGR to RGB color space to get hands on the video frame by frame as shown in the following code:

```
imgRGB = cv2.cvtColor(img, cv2.COLOR_BGR2RGB) self.results = self.hands.process(imgRGB) [8]
```


6.3 Mouse Control

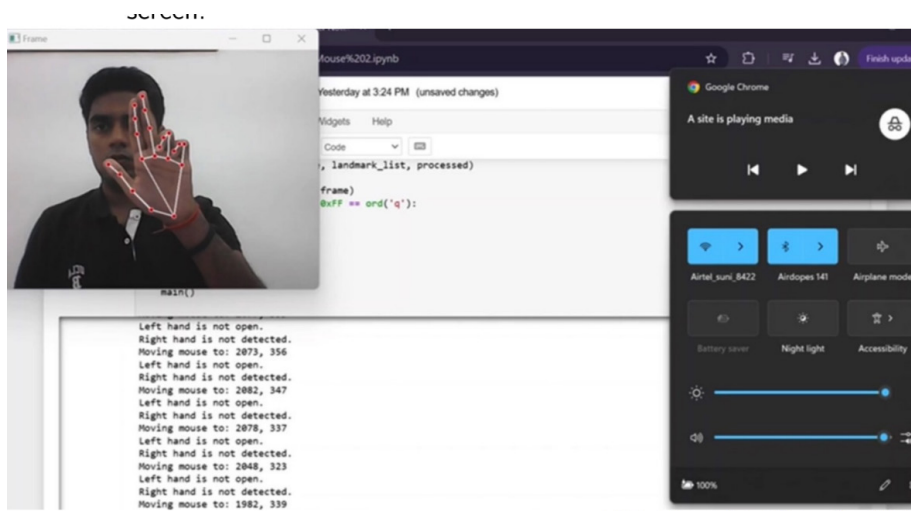


Figure 5: Mouse Control using Hand Gesture

Key Features of Mouse Control

- **Mouse Movement:** The ability to move the mouse pointer to a specific position on the screen (often represented by x and y coordinates). This is useful for automating GUI interactions where specific areas on the screen need to be clicked or hovered over. The position of the index finger tip and middle finger landmarks can be used to move the mouse cursor as shown in figure 5. By mapping the coordinates of the fingertip to the screen size, you can control the cursor's movement. To ensure smooth cursor movement, you can average out or interpolate the positions of the fingertip over multiple frames.
- **Mouse Click:** Simulating mouse button clicks (left, right, middle) at specific screen locations. This is essential for automating tasks that require user interaction, like clicking buttons, links, or menus. Index finger is used for 'left-click', middle finger is used for 'right-click' and both fingers are used for 'double-click'.

D. Gesture - based Custom Controls



Figure 6: Brightness Control using Hand Gesture

- **Brightness Control:** The ability to control brightness can be done by python libraries such as screen-brightness- control which can be used to control brightness. The screen-brightness-control library is a Python package that allows you to interact with and control the screen brightness on various operating systems, such as Windows, Linux, and MacOS. It provides a simple and easy-to-use interface to get the current screen brightness and adjust it programmatically.
- The system will detect if the both hands are raised. Once the both hands are identified as active, the gesture control logic will be activated. By measuring the distance between the index finger and thumb of the right hand, you can interpret this distance as a

way to adjust the brightness. The closer the fingers are, the lower the brightness, and vice versa.

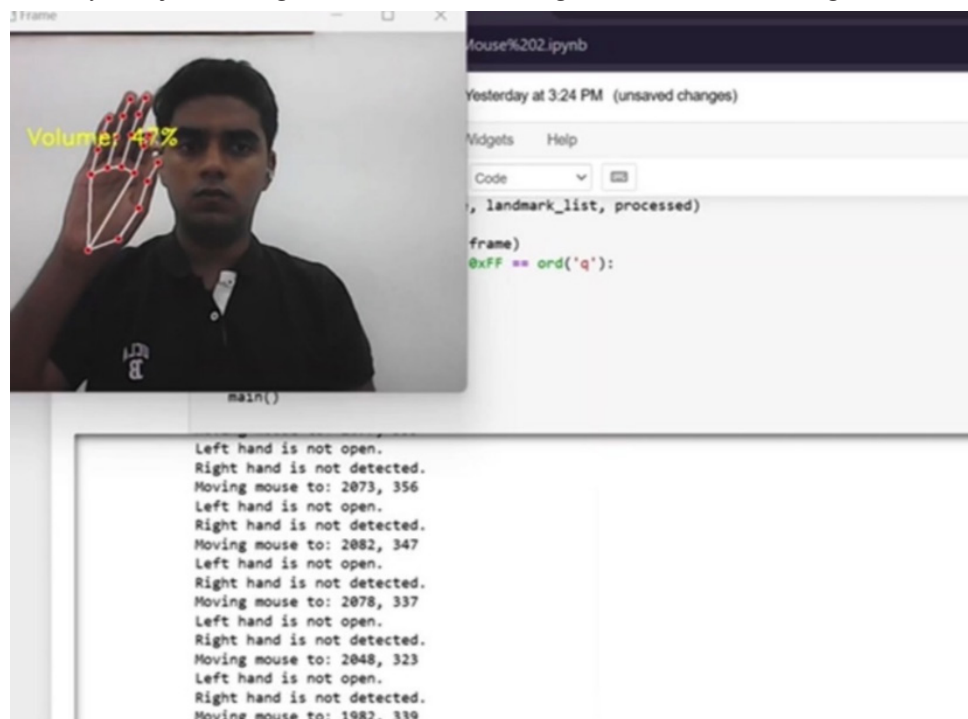


Figure 7: Volume Control using Hand Gesture

- **Volume Control:** The ability to control volume can be done by python libraries such as Pycaw which can be used to control volume. Pycaw is a Python library that provides an interface for controlling the audio properties of a Windows system. It utilizes the Windows Core Audio API to give developers access to various features, such as adjusting volume levels, retrieving device information, and managing audio sessions.
- The system will detect if the left hand is raised. Once the left hand is identified as active, the gesture control logic will be activated. By measuring the distance between the index finger and thumb of the left hand, you can interpret this distance as a way to adjust the volume. The closer the fingers are, the lower the volume, and vice versa.

7. Conclusion

- The visual AI mouse system revolutionizes user interaction by using a webcam to translate hand gestures into mouse functions like cursor movement, clicking, scrolling, and adjusting settings, eliminating traditional hardware and providing a more ergonomic and hygienic experience.
- This AI system shows high accuracy compared to existing models, making it ideal for minimizing physical contact and assisting those with mobility impairments, while also addressing hygiene concerns in post-pandemic environments.
- However, challenges remain, such as inconsistent right-click functionality and less precise actions for tasks like text selection or dragging, which can impact user efficiency and satisfaction.
- The system's potential lies in refining fingertip detection and improving accuracy for precision tasks, making it well-suited for environments prioritizing reduced contact.
- With further research, this touchless computing innovation could significantly enhance human-computer interaction, offering a seamless, accessible, and intuitive user experience.

References

1. Rao, A. K., & Gordon, A. M. (2001). Contribution of tactile information to accuracy in pointing movements. *Experimental brain research*, 138, 438-445.
2. Lira, M., Egito, J. H., Dall'Agnol, P. A., Amodio, D. M., Gonçalves, Ó. F., & Boggio, P. S. (2017). The influence of skin colour on the experience of ownership in the rubber hand illusion. *Scientific Reports*, 7(1), 1-13.
3. Jaloza, L. B. (2021). Inside Facebook Reality Labs: wrist-based interaction for the next computing platform. *Tech at Meta*, 18.
4. Carlton, B. (2021). HaptX Launches True-Contact Haptic Gloves For VR And Robotics. *VRScout*. URL <https://vrscout.com/news/haptx-truecontact-hapticgloves-vr/> (accessed 3.10. 21).
5. Brenton, H., Gillies, M., Ballin, D., Chatting, D. (2005). The uncanny valley: does it exist. In *Proceedings of conference of human computer interaction, workshop on human animated character interaction*. Pennsylvania: Citeseer.
6. Katona, J. (2021). A review of human-computer interaction and virtual reality research fields in cognitive InfoCommunications. *Applied Sciences*, 11(6), 2646.
7. Matlani, R., Dadlani, R., Dumbre, S., Mishra, S., Tewari, A. (2021). Virtual mouse using hand gestures. In *2021 international conference on technological advancements and innovations (ICTAI)* (340-345). IEEE.
8. Kumar, A., Chaudhary, A., Singhal, A., Dev, A., Jaiswal, A. (2024). Gesture and Voice Controlled Virtual Mouse: A Review. In *2024 2nd International Conference on Disruptive Technologies (ICDT)* (876-879). IEEE.
9. M. S. Inba Mary, M. Sakthi Anand, A. Sakthi Manikandan, and M. Senthamiluthu, (2024) "Hand gesture recognition using Mediapipe and OpenCV," *IJCRT*, 12 (2), 1-6.
10. Lugaresi, C., Tang, J., Nash, H., McClanahan, C., Uboweja, E., Hays, M., ... & Grundmann, M. (2019). Mediapipe: A framework for building perception pipelines. *arXiv preprint arXiv:1906.08172*.
11. Kumar, R., Bajpai, A., & Sinha, A. (2023). Mediapipe and cnns for real-time asl gesture recognition. *arXiv preprint arXiv:2305.05296*.
12. V. Bazarevsky and G. R. Fan Zhang. On-Device, Media Pipe for Real-Time Hand Tracking.
13. Tanner, T. (2019). Python GUI Automation Using PyAutoGUI. *Journal of Automation Scripting*.
14. Williams, P. (2018). Mathematical Techniques for GUI Automation. *Journal of Computer Science and Mathematics*.
15. Raffel, C., and Shinn, C. (2017). Human-Computer Interaction through Gestures. *IEEE Access*, 5, 4253-4263. DOI: 10.1109/AC-CESS.2017.2678859.
16. E. Sankar, B. N. Bharadwaj, and A. V. Vignesh. (2023). Virtual Mouse Using Hand Gesture. *International Journal of Scientific Research in Engineering and Management (IJSREM)*, 7(5), 1.

Copyright: ©2025 Aaditya Goel. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.