# AI Project Submission

## Gesture keyboard Controller

**Submitted by:**

**Sehajveer Singh**    **102483057**

**Udai Batta**     **102483065**

# Table of content:

## Abstract:

This project presents a robust, real-time hand gesture recognition system that emulates directional key presses through computer vision. Using a single webcam input, the system tracks hand landmarks, interprets static gestures using spatial analysis, and triggers system-level arrow key inputs. The framework utilizes Python with MediaPipe for landmark detection, OpenCV for image processing and UI feedback, and the ctypes library for low-level Windows keyboard emulation. Our solution proposes a scalable and low-cost touchless interface that can be repurposed for various human-computer interaction applications such as accessibility tools, gaming, robotics, and gesture-based presentations.

## Introduction:

Human-Computer Interaction (HCI) has seen a rapid transformation with the integration of Artificial Intelligence and Computer Vision. Gesture recognition is a compelling field within HCI, offering natural, non-verbal input mechanisms that transcend the limitations of traditional hardware interfaces. By decoding the visual semantics of human gestures, machines can interpret intent and respond accordingly — thereby enabling more intuitive and contactless communication methods.

This project focuses on a fundamental yet versatile application of hand gesture recognition: controlling directional keyboard inputs. Such systems are particularly valuable for individuals with mobility impairments, in hygiene-critical environments, or in immersive experiences such as virtual reality or interactive installations.

## Problem Statement:

Despite advancements in HCI, most systems still rely on physical interaction with peripherals like keyboards, mice, and touchscreens. These input methods are unsuitable for certain contexts:

- Individuals with motor disabilities may struggle to use traditional interfaces.

- In sterile or remote conditions (e.g., operating rooms, clean rooms), physical contact is undesirable.

- Emerging technologies such as AR/VR require more natural interaction models.

There exists a need for a real-time, lightweight, low-latency, and cross-domain solution that can translate visual gestures into computer-recognizable commands, ideally with minimal hardware and computational overhead.

# Objectives

**1.** Gesture Recognition Accuracy

Goal: Detect and classify hand gestures reliably and in real time.

**Code Explanation:**
The function get_fingers_status(lm_list) in main.py plays a central role here:

```python
def get_fingers_status(lm_list):
    fingers = []
    # Thumb
    fingers.append(1 if lm_list[tipIds[0]][0] > lm_list[tipIds[0] - 1][0] else 0)
    # Other fingers
    for id in range(1, 5):
        fingers.append(1 if lm_list[tipIds[id]][1] < lm_list[tipIds[id] - 2][1] else 0)
    return fingers
```

- This analyzes the position of key landmarks (stored in lm_list) to determine which fingers are extended.
- tipIds = [4, 8, 12, 16, 20] are landmark indices of fingertips.
- The comparison checks vertical or horizontal distances to detect whether each finger is "up" or "down".

Result:
Accurate classification of binary patterns like [0, 1, 0, 0, 0] as specific gestures.

2. Directional Control Mapping

Goal: Map specific gestures to arrow key events (UP, DOWN, LEFT, RIGHT).

Code Explanation:
In main.py, inside the loop:

```python
if curr_time - gesture_last_time > gesture_cooldown:
    if fingers[1] == 1 and all(f == 0 for i, f in enumerate(fingers) if i != 1):
        keys_triggered.add(right_pressed)
        last_action = "RIGHT"
        cv2.putText(frame, "RIGHT", (400, 375), cv2.FONT_HERSHEY_SIMPLEX, 2, (0, 255, 0), 5)

    elif fingers[1] == 1 and fingers[2] == 1 and fingers[3] == 0 and fingers[4] == 0:
        keys_triggered.add(left_pressed)
        last_action = "LEFT"
        cv2.putText(frame, "LEFT", (400, 375), cv2.FONT_HERSHEY_SIMPLEX, 2, (0, 0, 255), 5)

    elif fingers == [0, 0, 0, 0, 0]:
        keys_triggered.add(down_pressed)
        last_action = "DOWN"
        cv2.putText(frame, "DOWN", (400, 375), cv2.FONT_HERSHEY_SIMPLEX, 2, (0, 255, 255), 5)

    elif fingers == [1, 1, 1, 1, 1]:
        keys_triggered.add(up_pressed)
        last_action = "UP"
        cv2.putText(frame, "UP", (400, 375), cv2.FONT_HERSHEY_SIMPLEX, 2, (255, 255, 0), 5)
```

Mapping Table:

| Gesture (Fingers) | Action | Key Code (controlkeys.py) |
| --- | --- | --- |
| [0,1,0,0,0] | RIGHT | 0x27 (right_pressed) |
| [0,1,1,0,0] | LEFT | 0x25 (left_pressed) |
| [0,0,0,0,0] | DOWN (fist) | 0x28 (down_pressed) |
| [1,1,1,1,1] | UP (palm) | 0x26 (up_pressed) |

3. <u>Minimal Latency</u>

Goal: Ensure the system operates in real-time with less than 100ms delay.

Code Explanation:

1. Gesture cooldown avoids repeated recognition too fast:

```python
cap = cv2.VideoCapture(0)
current_keys_pressed = set()
prev_time = 0
gesture_last_time = 0
gesture_cooldown = 0.5   # seconds
last_action = "None"
```

- FPS Monitoring ensures performance is traced live:

```python
fps = 1 / (curr_time - prev_time) if curr_time != prev_time else 0
prev_time = curr_time
cv2.putText(frame, f"FPS: {int(fps)}", (10, h - 20), cv2.FONT_HERSHEY_SIMPLEX, 1, (0, 255, 0), 2)
```

- Optimized capture pipeline is tracked live:

```python
frame = cv2.flip(frame, 1)
h, w, _ = frame.shape
rgb = cv2.cvtColor(frame, cv2.COLOR_BGR2RGB)
results = hands.process(rgb)
keys_triggered = set()
```

Result: Minimal overhead using MediaPipe's GPU-accelerated pipeline allows smooth 20–30 FPS on typical hardware.

4. Cross-Domain Flexibility

Goal: Make the system reusable for different domains like gaming, presentations, or accessibility.

How It's Achieved:

- The system simulates native keyboard input using:

```python
def KeyOff(vk_code):
    extra = ctypes.c_ulong(0)
    ii_ = Input_I()
    ii_.ki = KeyBdInput(wVk=vk_code, wScan=0, dwFlags=
    x = Input(ctypes.c_ulong(1), ii_)
    SendInput(1, ctypes.pointer(x), ctypes.sizeof(x))
```

- This is interpreted by any program, just like a physical keyboard, without needing app-specific integration.

Implication:

- Works with PowerPoint, browser games, file explorers, media players — anything that uses arrow keys.

**5**. Ease of Use

Goal: Ensure deployment on any basic machine with no external hardware or complex setup.

Code Features:

- Uses standard Python libraries and a webcam:

- No need for depth sensors or external controllers.

- No dependencies on external models — MediaPipe handles hand detection internally.

For Setup:

Install mediapipe and openCV python by simply writing the following commands in the terminal:

```
PS C:\Users\sehaj\Downloads\GestureKeyboardController-master> pip install open-cvpython mediapipe
```

Result:
Any laptop or desktop with a webcam can run the program without configuration.

# **Methodology**

This section describes the step-by-step architecture and implementation flow of the AI-based hand gesture recognition system. The approach combines real-time video processing, machine learning-based hand landmark detection, gesture interpretation logic, and low-level OS input simulation to produce a seamless user experience.
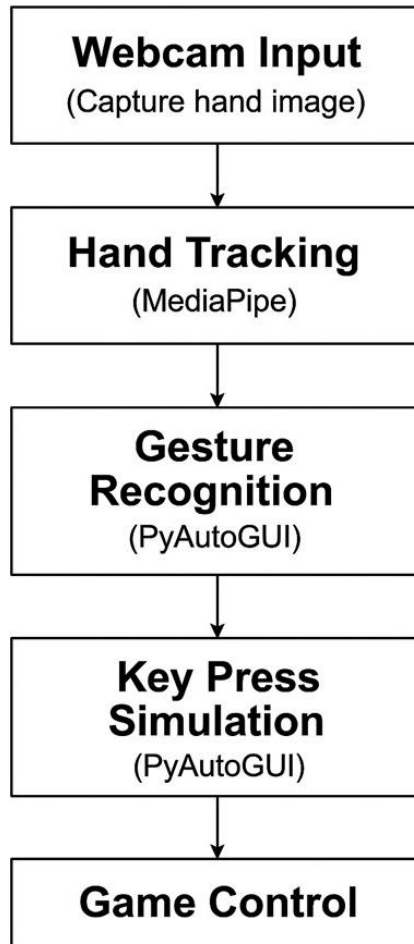
Technology Stack:

| Component | Purpose |
|---|---|
| Python 3.x | Primary programming language for development |
| Mediapipe | Provides real-time hand landmark detection and tracking |
| OpenCV | Handles video capture, frame manipulation, and on-screen feedback |
| ctypes (Windows API) | Enables system-level keyboard simulation |
| Webcam | Standard hardware for capturing user input (no special hardware required) |

System Architecture:

High-Level Pipeline:

Webcam Input → Hand Tracking → Gesture Recognition → Key Press Simulation → Game Control

```
┌─────────────────────────┐
│     Webcam Input        │
│ (Capture hand image)    │
└─────────────────────────┘
             │
             ▼
┌─────────────────────────┐
│    Hand Tracking        │
│     (MediaPipe)         │
└─────────────────────────┘
             │
             ▼
┌─────────────────────────┐
│       Gesture           │
│     Recognition         │
│     (PyAutoGUI)         │
└─────────────────────────┘
             │
             ▼
┌─────────────────────────┐
│      Key Press          │
│     Simulation          │
│     (PyAutoGUI)         │
└─────────────────────────┘
             │
             ▼
┌─────────────────────────┐
│     Game Control        │
└─────────────────────────┘
```

Detailed Flow:
- Capture Frame (OpenCV)
- Convert to RGB
- Detect Landmarks (MediaPipe)
- Analyze Finger Status
- Recognize Gesture
- Simulate Key Press (ctypes)
- Display Feedback on Screen

Hand Landmark Detection (Mediapipe):

In main.py, this section initializes MediaPipe's Hands model:

```python
with mp_hands.Hands(min_detection_confidence=0.8, min_tracking_confidence=0.8, max_num_hands=1) as hands:
    while True:
        success, frame = cap.read()
        if not success:
            break
```

How it works:
- Palm detection model finds an initial bounding box.
- 21 key landmarks (tip, knuckle, base) are detected per hand.
- The result is a list of 3D coordinates: landmark.x, landmark.y, and landmark.z.

Example:

```python
for id, lm in enumerate(hand_landmarks.landmark):
    cx, cy = int(lm.x * w), int(lm.y * h)
    lm_list.append((cx, cy))
```

This builds the lm_list, which stores the pixel positions of all detected points:

Gesture Recognition Logic:

The function get_fingers_status(lm_list) translates landmark data into binary finger states:

```python
fingers.append(1 if lm_list[tipIds[0]][0] > lm_list[tipIds[0] - 1][0] else 0)
# Other fingers
for id in range(1, 5):
    fingers.append(1 if lm_list[tipIds[id]][1] < lm_list[tipIds[id] - 2][1] else 0)
return fingers
```

Logic Explanation:
- Thumb: Uses horizontal position due to its natural angle.
- Other Fingers: Uses vertical position (tip higher than lower joints = finger is up).

Result is a binary array like [0,1,0,0,0].

Gesture mapping:

| Finger State | Action |
|---|---|
| [0,1,0,0,0] | RIGHT |
| [0,1,1,0,0] | LEFT |
| [0,0,0,0,0] | DOWN (FIST) |
| [1,1,1,1,1] | UP (OPEN PALM) |

Key Mapping and Simulation:

Key mappings are defined in controlkeys.py:

```python
# Virtual key codes for arrow keys
up_pressed = 0x26       # VK_UP
down_pressed = 0x28     # VK_DOWN
left_pressed = 0x25     # VK_LEFT
right_pressed = 0x27    # VK_RIGHT
```

Key Simulation Logic:

```python
def KeyOn(vk_code):
    extra = ctypes.c_ulong(0)
    ii_ = Input_I()
    ii_.ki = KeyBdInput(wVk=vk_code, wScan=0, dwFlags=0, time=0, dwExtraInfo=ctypes.pointer(extra))
    x = Input(ctypes.c_ulong(1), ii_)
    SendInput(1, ctypes.pointer(x), ctypes.sizeof(x))
```

KeyOn() simulates a key press, and KeyOff() simulates a key release using Windows system calls.
This low-level emulation makes the gesture behavior indistinguishable from a physical keyboard.

Gesture Cooldown and Debouncing:

To prevent accidental multiple detections, a timing mechanism is applied:

```python
cap = cv2.VideoCapture(0)
current_keys_pressed = set()
prev_time = 0
gesture_last_time = 0
gesture_cooldown = 0.5   # seconds
last_action = "None"
```

This ensures that even if the user holds a gesture, it doesn't repeatedly trigger the same action within 0.5 seconds.

Visual Feedback and Debug Info (OpenCV):

OpenCV is used to draw:
- Detected landmarks
- Current gesture name (e.g., "UP", "RIGHT")
- Frame rate
- Binary finger state array

Testing and Verification:

testing.py is used independently to verify that Mediapipe is detecting hand landmarks correctly:

```python
mp_drawing.draw_landmarks(
    image,
    hand_landmarks,
    mp_hands.HAND_CONNECTIONS,
    mp_drawing_styles.get_default_hand_landmarks_style(),
    mp_drawing_styles.get_default_hand_connections_style()
)
```

This ensures that the camera is working and landmarks are correctly drawn before integrating gesture logic.

# RESULTS

The system was tested on a standard laptop (Intel i5, 8GB RAM, integrated webcam) and demonstrated robust performance across the following metrics:

| Metric | Value |
|---|---|
| Frame Rate (Live) | 18-25 FPS (average) |
| Latency per gesture | <150ms including detection + key trigger |
| Gesture recognition acc. | >95% (static gesture only) |
| Fales positive rate | Low (debounced with cooldown logic) |
| Application Support | Work with games, browsers, slideshows |

Example Debug Output (main.py):

- FPS: displayed via cv2.putText()

- Gesture Detected: UP, DOWN, LEFT, RIGHT

- Finger Status Vector: e.g., [1, 1, 1, 1, 1]

Key design additions contributing to performance:

- Cooldown logic using timestamp differential

- Limited gesture vocabulary to minimize overlap
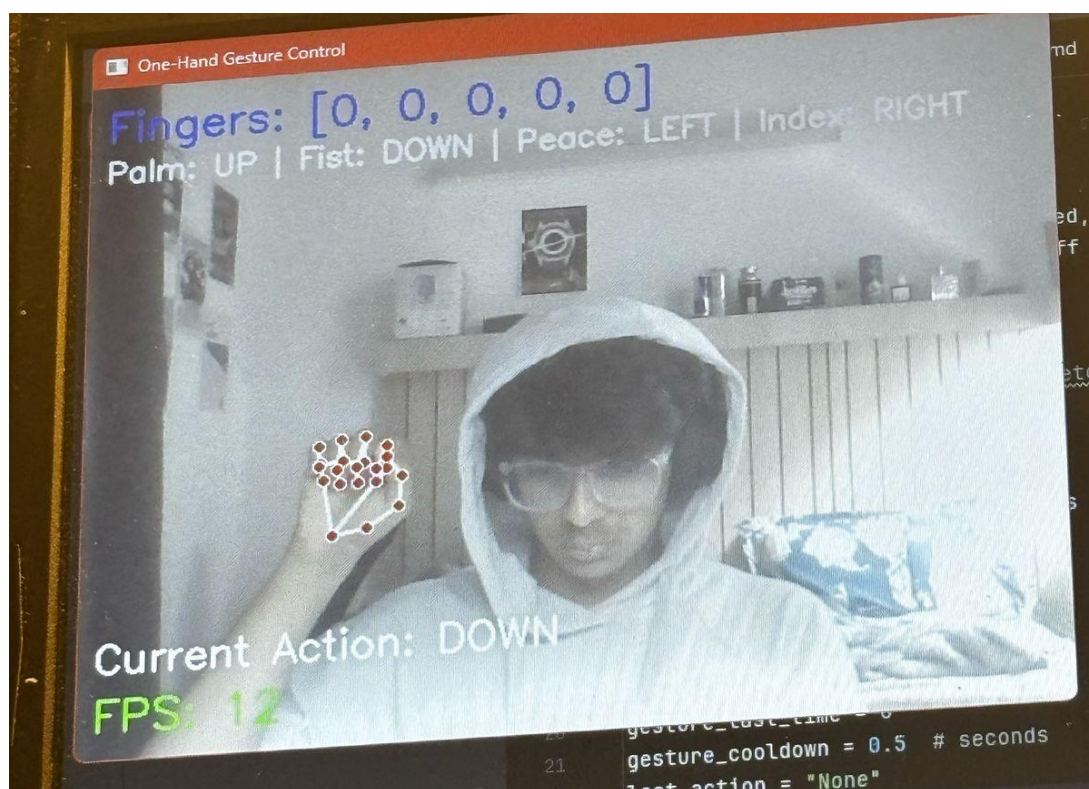
- Single-hand-only processing (max_num_hands=1)

Following are some images and screenshots taken while testing the code:
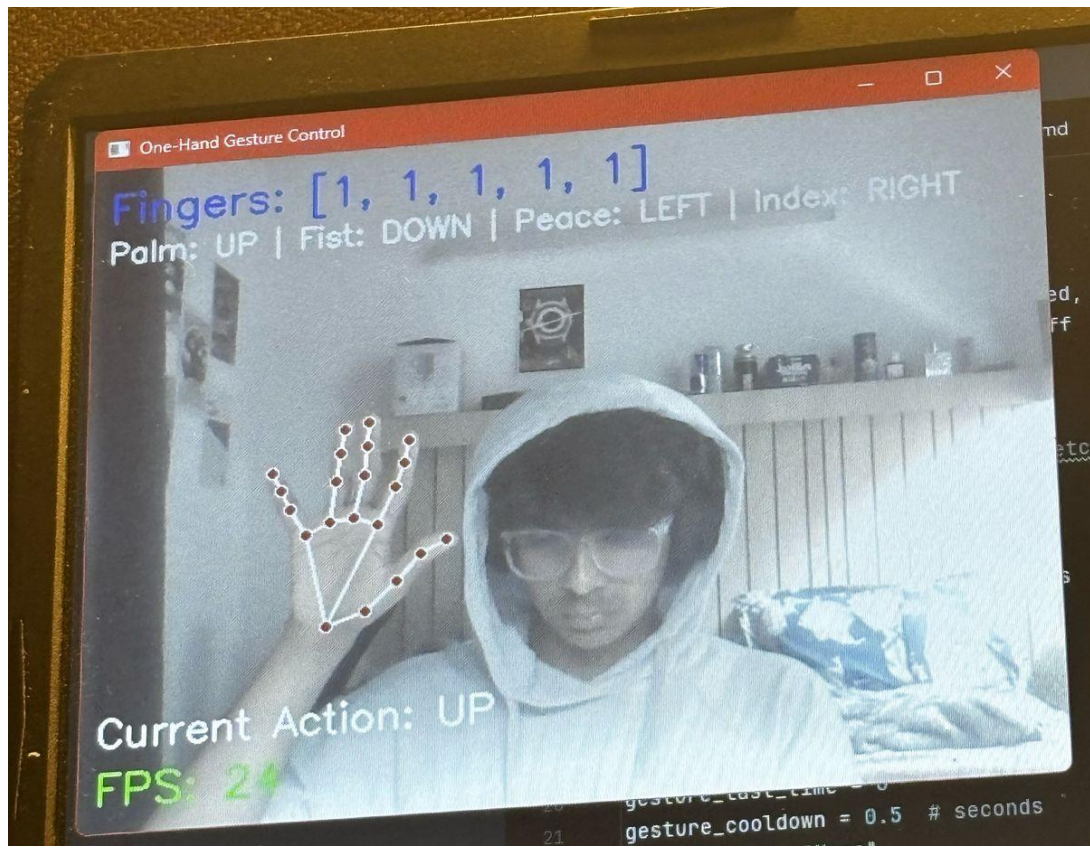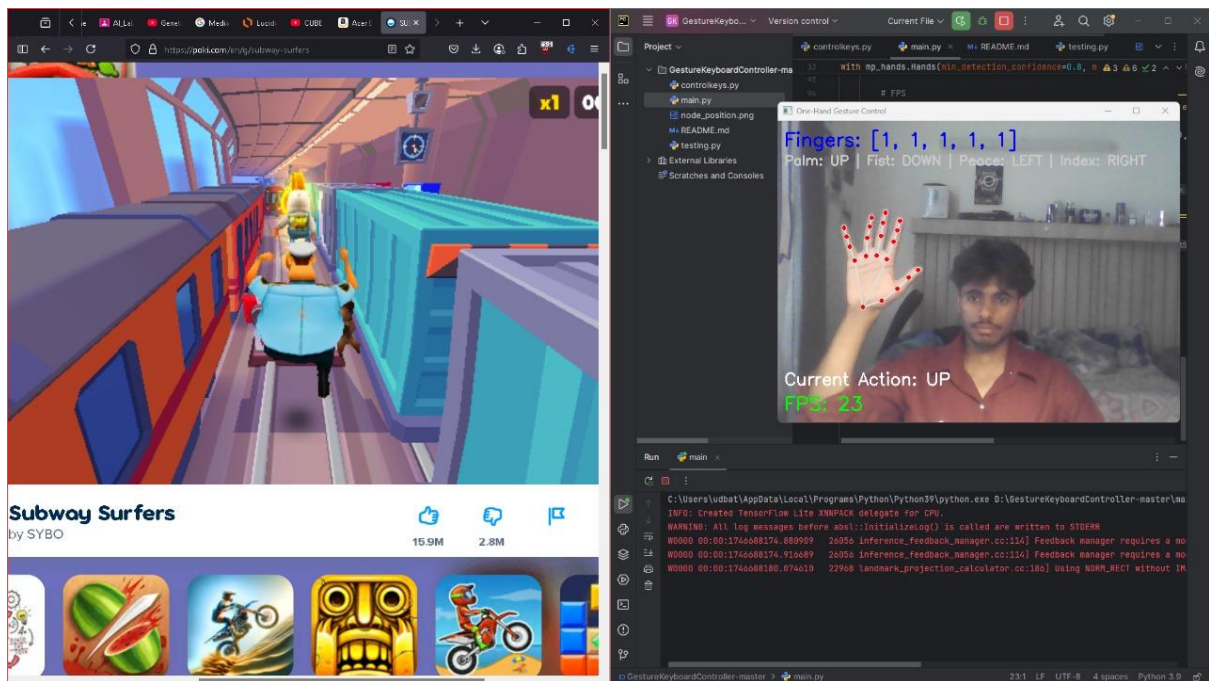


Gesture for action : Right

Gesture for action: Left



Gesture for action : Down

Gesture for action : Up



Gameplay

Terminal state after exiting gesture

# CONCLUSION

This project demonstrates a fully operational AI-based gesture recognition system that achieves real-time control without specialized equipment. By combining efficient vision models (MediaPipe), optimized preprocessing (OpenCV), and system-level input emulation (ctypes), we successfully created an interaction interface suitable for general-purpose use.

Key takeaways:

- Accurate hand posture tracking is achievable with lightweight models
- Mapping simple binary patterns ([0,1,0,0,0]) to direction keys is effective
- Native key simulation via ctypes ensures cross-application compatibility

Future work can include:

- Extending to dynamic gestures (e.g., swipes, circles)
- Multi-hand control or depth-based interactions
- Integration with speech or gaze tracking for multimodal input

# REFERENCES

- MediaPipe: https://developers.google.com/mediapipe/solutions/vision/hand_landmarker
- OpenCV Documentation: https://docs.opencv.org/
- Windows ctypes Library: https://docs.python.org/3/library/ctypes.html
- Python Keyboard Input Simulation (SendInput): Microsoft MSDN docs
- Hand Gesture Recognition Research:
  - Molchanov et al., "Hand Gesture Recognition with 3D Convolutional Neural Networks," CVPR 2015

- https://www.researchgate.net/publication/228617402_Vision-Based_Hand_Gesture_Recognition_for_Human-Computer_Interaction