

Mouse control using eyegaze

Aryaman Oberoi
IIT Kanpur
aryamano24@iitk.ac.in

Chandraveer Sisodia
IIT Kanpur
cssisodia24@iitk.ac.in

Shivam Kumar
IIT Kanpur
shivamkk24@iitk.ac.in

Anikait Dixit
IIT Kanpur
anikaitd24@iitk.ac.in

Abstract—This study develops a robust method for hands-free mouse cursor control, addressing the unreliability of direct eye-gaze tracking on common webcams. We first implement and compare two distinct frameworks, Dlib and MediaPipe, to track facial landmarks and calculate head pose using a Perspective-n-Point (PnP) algorithm. Our comparative analysis demonstrates that MediaPipe is objectively superior, offering more than double the Frames Per Second (FPS) and overcoming Dlib’s critical failure to track non-frontal faces.

Based on these conclusive findings, we then constructed a final, optimized controller built entirely on the MediaPipe framework. This novel system combines stable MediaPipe Face Mesh (for head-pose navigation) with MediaPipe Hands to detect both fist gestures (for reliable clicks) and pinch gestures (for scrolling). This unified, high-performance approach results in a stable, low-latency, and practical hands-free interface, making it a superior solution for accessible human-computer interaction and assisting individuals with severe Motor Disabilities, Spinal Cord Injuries, or Neurodegenerative diseases.

I. INTRODUCTION

In recent years, advancements in human-computer interaction (HCI) have led to the development of alternative input methods beyond traditional keyboard and mouse interfaces. Among these, intuitive, hands-free solutions are emerging with immense potential, particularly for individuals with motor impairments, offering them a more accessible way to navigate and interact with digital systems.

Our initial exploration focused on direct eye gaze tracking (which follows the pupil), as it is a powerful concept. However, we immediately encountered significant real-world failures. Our attempts to implement direct iris tracking using a standard, integrated laptop webcam were met with extreme instability. The poor quality and low resolution of the video feed made it impossible for the algorithms to maintain a stable lock on the pupil’s small, noise-sensitive movements. This resulted in an unusable, “jittery” cursor that would jump erratically across the screen.

Similarly, our efforts to use blink detection for clicks were unsuccessful for the same reason. The low-fidelity video made it difficult to reliably distinguish an intentional blink from an open eye, leading to a frustrating experience with constant false positives (accidental clicks) and missed inputs.

Faced with these practical limitations, this project directly addresses this challenge by pivoting to a more robust alternative: head pose estimation. Instead of tracking the unreliable movements of the iris, our system analyzes the head’s overall orientation (specifically, its yaw and pitch) as a stable proxy for the user’s gaze direction. This approach is far more

resilient and functions effectively even with the common, low-resolution webcams that made eye tracking impossible.

A primary challenge in developing such a real-time system is achieving an optimal balance between precision (is the cursor accurate enough to be useful?) and computational efficiency (does it run smoothly without lagging?). To find this balance, our research implements and compares two distinct computer vision frameworks for the core head-pose task: the classic, regression-tree-based Dlib library and the modern, deep-learning-based MediaPipe framework.

Furthermore, to create a complete and functional interface, we augment this head-pose navigation with robust gesture controls. We replace the failed traditional blink detection by integrating MediaPipe Hands into our system. This module is used to detect two specific gestures: a fist gesture for reliable clicks and a pinch gesture for scrolling. Our goal is to analyze this combined-model approach to develop an accessible, efficient, and user-friendly alternative to conventional input devices.

II. LANDMARK DETECTION

In computer science, landmark detection is the process of finding significant landmarks in an image. This originally referred to finding landmarks for navigational purposes – for instance, in robot vision or creating maps from satellite images. Methods used in navigation have been extended to other fields, notably in facial recognition where it is used to identify key points on a face. It also has important applications in medicine, identifying anatomical landmarks in medical images.

A. Facial Landmark

Detecting facial landmarks is a subset of the shape prediction problem. Given an input image (and normally an ROI that specifies the object of interest), a shape predictor attempts to localize key points of interest along the shape.

In the context of facial landmarks, our goal is detect important facial structures on the face using shape prediction methods. Some of the most popular frameworks used are Dlib, Active Shape Models (ASM), and the more recent deep learning-based MediaPipe FaceMesh. Detecting facial landmarks is therefore a two step process:

- Step #1: Localize the face in the image.
- Step #2: Detect the key facial structures on the face ROI.

B. Dlib

The pre-trained facial landmark detector inside the dlib library is used to estimate the location of 68 (x, y)-coordinates that map to facial structures on the face.



Fig. 1. Dlib 68 facial landmarks visualization.

For head pose estimation, a subset of these 68 landmarks is selected, corresponding to distinct anatomical features that provide a stable reference for estimating orientation. The selected landmarks and their Dlib indices are:

- **Nose Tip (index 30):** Used as the central reference point of the face.
- **Left Eye Outer Corner (index 36) & Right Eye Outer Corner (index 45):** These points help measure the horizontal direction of the face (left–right alignment).
- **Left Mouth Corner (index 48) & Right Mouth Corner (index 54):** Useful for estimating yaw — whether the face is turned left or right.
- **Chin (index 8):** Helps with vertical orientation and pitch — whether the face is looking up or down.

1) *2D Coordinates of These Landmarks:* The 2D position of each selected landmark in the image is written as:

$$p_i = (x_i, y_i), \quad \text{for } i \in \{30, 36, 45, 48, 54, 8\} \quad (1)$$

Which simply means:

- Each landmark has an (x, y) coordinate,
- And we are using these six landmark indices: 30, 36, 45, 48, 54, 8.

These landmarks provide a stable geometric foundation for head pose estimation. This allows for the calculation of rotation angles (yaw, pitch, and roll) by contrasting their relative positions with a predefined 3D head model.

III. ARCHITECTURE OF MEDIAPIPE FACEMESH

A. Overview of MediaPipe

MediaPipe is an open-source framework that enables developers to construct high-performance machine learning (ML)

perception pipelines. The framework includes a library of adaptable modules for specific tasks, such as hand tracking, face detection, and pose estimation. These components are optimized for real-time operation on diverse platforms, from web browsers to mobile devices.

B. FaceMesh Architecture

MediaPipe’s FaceMesh module utilizes a lean convolutional neural network (CNN) to identify and map 468 distinct facial landmarks, all in real-time. This architecture is built on several key components:

- **Landmark Detection Network:** This network analyzes an input image to find facial features. It notably uses depthwise separable convolutions, a technique that lowers the computational cost without sacrificing accuracy. The model is also robust enough to work reliably under various conditions, such as different head orientations or facial expressions.
- **Facial Landmark Representation:** The 468 points provide a detailed map of the face, with each point corresponding to a specific feature like the eyes, nose, or mouth. Such a dense and comprehensive map enables in-depth analysis, like assessing facial symmetry.
- **Modular Pipeline Structure:** The framework is built modularly, which allows for the easy combination of different processing stages (e.g., face detection followed by landmark extraction). This pipeline design is engineered for highly efficient data handling, resulting in fast inference speeds and low latency—a crucial requirement for any real-time application.

C. Visualization Capabilities

The FaceMesh package also comes with integrated visualization tools. These utilities can be used to overlay the detected landmarks on an image, drawing features like the contours of the face, a full facial mesh (tessellation), or the iris, which provides clear and intuitive visual feedback.

D. Performance Optimizations

FaceMesh is specifically engineered for high performance on mobile and edge computing devices. This is achieved through techniques like model quantization and options for hardware acceleration. These strategies enable the system to run at high frame rates while using minimal computational resources, making it an ideal choice for demanding real-time tasks like augmented reality or live video conferencing.

The selected landmarks and their respective roles are:

- **Nose Tip (i = 1):** Acts as a central reference point for pose estimation.
- **Left Eye Outer Corner (i = 133) and Right Eye Outer Corner (i = 362):** Help determine horizontal alignment and yaw angle.
- **Left Ear (i = 234) and Right Ear (i = 454):** Provide additional spatial constraints to estimate head rotation.
- **Chin (i = 152):** Helps in estimating vertical orientation and pitch angle.

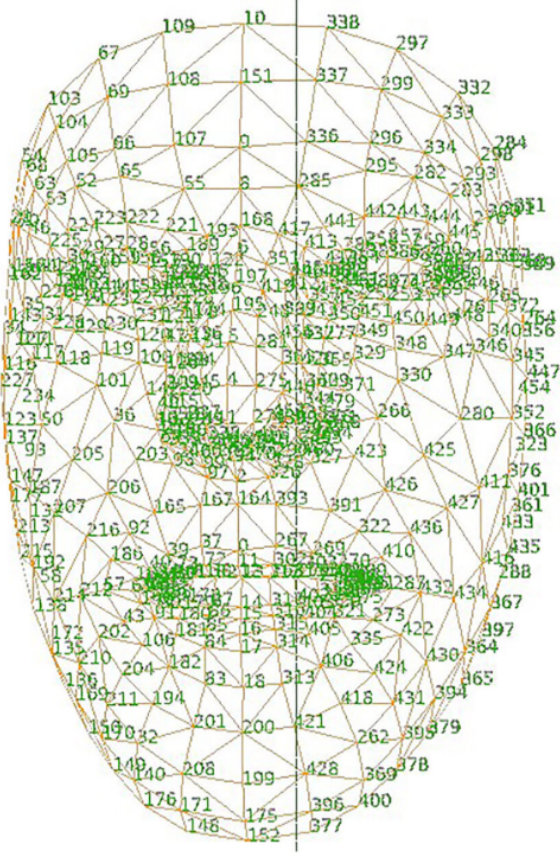


Fig. 2. MediaPipe FaceMesh 468 facial landmarks visualization.

The 2D image coordinates of these selected landmarks are given by:

$$p_i = (x_i, y_i), \quad i \in \{1, 133, 362, 234, 454, 152\} \quad (2)$$

These landmarks provide a stable geometric foundation for head pose estimation. This allows for the calculation of rotation angles (yaw, pitch, and roll) by contrasting their relative positions with a predefined 3D head model.

IV. METHODOLOGY

A. What is Pose Estimation?

In computer vision the pose of an object refers to its relative orientation and position with respect to a camera. You can change the pose by either moving the object with respect to the camera, or the camera with respect to the object.

The pose estimation problem described in this tutorial is often referred to as Perspective-n-Point problem or PNP in computer vision jargon. As we shall see in the following sections in more detail, in this problem the goal is to find the pose of an object when we have a calibrated camera, and we know the locations of n 3D points on the object and the corresponding 2D projections in the image.

B. Camera Calibration and Intrinsic Parameters

Accurate head pose estimation requires defining the camera's intrinsic matrix, K . This matrix mathematically models how 3D real-world points are projected onto the 2D image plane, assuming a standard pinhole camera model.

The intrinsic matrix K is defined as:

$$K = \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} \quad (3)$$

where:

- f_x and f_y are the focal lengths in pixels along the x and y axes
- c_x and c_y are the coordinates of the principal point (optical center)
- The 3D-to-2D scaling of the scene is governed by the pixel-based focal lengths, f_x and f_y . In our practical application, we employ a common simplification, setting both f_x and f_y to a value equivalent to the image's total width.
- Concurrently, the camera's optical axis, or principal point (c_x, c_y), is assumed to be at the precise geometric center of the video frame.

1) *Lens Distortion Assumption:* Physical camera lenses invariably introduce image warping, which manifests as radial (k_1, k_2, k_3) and tangential (p_1, p_2) distortions that visibly curve straight lines.

The distortion coefficient vector is:

$$D = [k_1, k_2, p_1, p_2, k_3] \quad (4)$$

Correcting this warp is computationally intensive. To prioritize real-time performance and simplicity, our project bypasses this by adopting an ideal pinhole camera model, which is theoretically free of all distortion. This assumption is implemented by setting the distortion coefficient vector D to a zero vector:

$$D = [0, 0, 0, 0, 0] \quad (5)$$

The direct benefit of this simplification is that the projection of 3D points onto the 2D image plane becomes a purely linear transformation dictated by the K matrix. This significantly streamlines the calculations required for pose estimation.

2) *Effect on Head Pose Estimation:* These intrinsic parameters are the bedrock of the solvePnP algorithm. The algorithm is critically dependent on the K matrix; it uses K as the "key" to translate the 3D model's points from their fixed world coordinates to the 2D pixel locations detected in the camera's feed.

While using the image width as a substitute for the true focal length is a pragmatic shortcut, it provides a "reasonable approximation" of the camera's perspective projection. For our goal of stable, real-time control, this approximation is both sufficient and necessary to allow the algorithm to derive an accurate rotation estimation.

C. Head Pose Estimation

To compute the head's orientation, we first require a generic 3D model of a face. This model is defined by a set of predefined 3D coordinates (P_i). These 3D points are strategically chosen based on their anatomical significance to provide a stable and accurate foundation for the pose estimation.

In our project, we implement two distinct 3D models (defined in our code as MODEL_POINTS) to align with the different landmark predictors:

- **For MediaPipe:** The 3D model is a 6-point array representing the Nose Tip, Left Eye, Right Eye, Left Ear, Right Ear, and Chin.
- **For Dlib:** The 3D model is a 6-point array representing the Nose Tip, Left Eye Outer Corner, Right Eye Outer Corner, Left Mouth Corner, Right Mouth Corner, and Chin.

These 3D model points (P_i) must correspond to the 2D pixel coordinates (p_i) of the facial landmarks detected in the live video feed. The specific landmark indices used to find these 2D points differ between our two frameworks:

- **MediaPipe 2D Indices:** $i \in \{1, 133, 362, 234, 454, 152\}$
- **Dlib 2D Indices:** $i \in \{30, 36, 45, 48, 54, 8\}$

With both sets of points—the 2D detected points (p_i) and their corresponding 3D model points (P_i)—we can calculate the head's orientation using the Perspective-n-Point (PnP) algorithm. The goal of the PnP algorithm is to find a solution for the main projection equation:

$$p_i = K \cdot (RP_i + t) \quad (6)$$

which can be expanded as:

$$\begin{bmatrix} u_i \\ v_i \\ 1 \end{bmatrix} \sim \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} \left(\begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix} \begin{bmatrix} X_i \\ Y_i \\ Z_i \end{bmatrix} + \begin{bmatrix} t_x \\ t_y \\ t_z \end{bmatrix} \right) \quad (7)$$

Where:

- K is the camera's intrinsic matrix, which was defined previously.
- R is the 3×3 rotation matrix. This is the primary output we need, as it mathematically describes the head's 3D orientation.
- t is the 3×1 translation vector, which accounts for the head's 3D position relative to the camera.
- To solve for R and t , we use the OpenCV solvePnP function, which employs methods such as the Levenberg-Marquardt algorithm or RANSAC-based iterative refinement.

D. Conversion to Euler Angles

The rotation matrix R , which is the output of the PnP algorithm, mathematically describes the head's 3D orientation. However, this 3×3 matrix is not directly interpretable for our control purposes. To make this orientation data human-readable, we must convert the matrix R into the more intuitive Euler angles.

The rotation matrix R has the form:

$$R = \begin{bmatrix} R_{11} & R_{12} & R_{13} \\ R_{21} & R_{22} & R_{23} \\ R_{31} & R_{32} & R_{33} \end{bmatrix} \quad (8)$$

We use the following formulas to extract the yaw, pitch, and roll:

$$\text{yaw} = \arctan 2(R_{21}, R_{11}) \quad (9)$$

$$\text{pitch} = \arctan 2(-R_{31}, \sqrt{R_{32}^2 + R_{33}^2}) \quad (10)$$

$$\text{roll} = \arctan 2(R_{32}, R_{33}) \quad (11)$$

These angles are defined as:

- **Yaw (θ_{yaw}):** Rotation around the vertical axis, corresponding to a left-right head movement.
- **Pitch (θ_{pitch}):** Rotation around the horizontal axis, corresponding to an up-down head movement.
- **Roll (θ_{roll}):** Rotation around the depth axis, corresponding to tilting the head sideways.

E. Smoothing of Values

The raw yaw and pitch angles computed from the PnP algorithm can be "jittery" due to small, noisy variations in the detected landmarks from frame to frame. To fix this and create a stable cursor, we apply an exponential moving average (EMA) to smooth the computed angles.

This is implemented in our code for both smoothed_yaw and smoothed_pitch. The formula is:

$$\theta^{(t)} = \alpha \theta_{\text{new}} + (1 - \alpha) \theta^{(t-1)} \quad (12)$$

where:

- $\theta^{(t)}$ is the new smoothed angle.
- $\theta^{(t-1)}$ is the previously smoothed angle.
- θ_{new} is the new, raw angle calculated from the PnP algorithm.
- α is the smoothing parameter.

In our implementation, we set $\alpha = 0.7$. This α value represents the "weight" given to the new, raw data.

- A higher alpha (like our 0.7) gives more weight to the new values, resulting in a more responsive but slightly "noisier" cursor.
- A lower alpha would give more weight to the previous smoothed values, resulting in a much smoother but slower, "lagging" response.

Our value of 0.7 was chosen as a good balance between responsiveness and stability for controlling the mouse.

F. Gaze Visualization

The estimated gaze direction is visualized as a vector \mathbf{g} projected from the face center. Given the smoothed yaw and pitch angles, we approximate the gaze direction as:

$$\mathbf{g} = (g_x, g_y) = (-L \sin \theta_{\text{yaw}}, -L \sin \theta_{\text{pitch}}) \quad (13)$$

where:

- L is the predefined arrow length (typically a fixed scalar).
- θ_{yaw} and θ_{pitch} control the horizontal and vertical displacement of the arrow, respectively.

The endpoint of the gaze vector is computed as:

$$p_{\text{end}} = p_{\text{center}} + \mathbf{g} \quad (14)$$

where p_{center} is the face center, which is often approximated as the nose tip landmark (p_1). This allows real-time visualization of where the person is looking within the image frame.

1) *Significance of Gaze Estimation:* Gaze estimation is a foundational technology for a new class of interactive applications. The most critical applications for our project include:

- **Assistive Technologies:** This is the primary motivation. It provides a crucial "hands-free" interface, enabling individuals with significant motor impairments to access and control computers.
- **Human-Computer Interaction (HCI):** Our project is a core example of this, using "gaze-based controls" to create a new and intuitive way for users to interact with their systems.
- **Immersive Gaming and Simulation:** The same head-pose tracking can be used to control the in-game camera (e.g., in a flight simulator) or interact with virtual environments, offering a more engaging experience.

Our project's methodology, which focuses on PnP and smoothing, is a direct effort to achieve "stable and real-time gaze tracking," which is essential for making any of these applications robust and user-friendly.

V. COMPARISON BETWEEN THE TWO MODELS

We have compared the two frameworks side by side quantitatively using FPS for both the models.

A. FPS

Frames Per Second (FPS) is a performance metric that measures how many individual images, or "frames," a system can process and render in a single second.

1) *Why FPS is a Good Metric for Comparison:* FPS is an excellent metric for this project because it directly quantifies the computational efficiency and user-perceived responsiveness of each framework. As High FPS means low latency and fluid cursor motion.

B. Dlib Framework

Observations - Average FPS: 14-30 frames per second.

1) *Dlib Pose-Tracking Failure:* A major limitation was observed in the Dlib implementation: the cursor tracking completely stops when the head is turned even slightly. `dlib.get_frontal_face_detector()` is only trained to find "frontal" faces. The moment you turn your head, your face becomes a "profile" view, and the detector fails to find a face.

C. MediaPipe Framework

Observations - Average FPS: 100-135 frames per second.

The tracking failure for high values of Yaw and pitch were not observed, as the MediaPipe FaceMesh model is designed to be robust to varying head poses and does not suffer from this "frontal-only" limitation.



Fig. 3. Dlib framework performance showing low FPS.

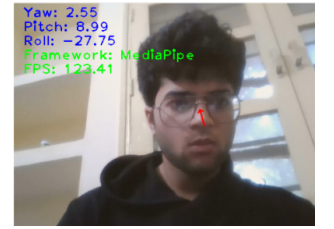


Fig. 4. MediaPipe framework performance showing high FPS.

TABLE I
COMPARISON OF DLIB AND MEDIAPIPE

Feature	Dlib (Facial Landmark Detector)	MediaPipe (Face Mesh)
Core Technology	"Classic" Machine Learning: HOG + Regression Trees	Modern Deep Learning (Lightweight CNN)
Performance (FPS)	CPU-heavy, not optimized for real-time. Runs 15-30 FPS, which feels laggy	Highly optimized for real-time. Easily runs at 110+ FPS
Landmarks	Provides 68 2D landmarks (x, y coordinates)	Provides 478 3D landmarks (x, y, and z coordinates)
Robustness	Trained on frontal faces. Fails when head turned too far	Robust to head pose, lighting, and partial blockages
Model Files	External. Requires manual download of 80MB+ .dat file	Self-contained. Model bundled directly into library

VI. CLICK TRIGGER

Our system uses MediaPipe Hands, a high-fidelity, real-time hand-tracking model that provides 21 3D landmarks (or keypoints) for each hand, to serve as the "event detection output". We then use the `pyautogui` Python library for programmatic "control of the mouse";

upon detecting a valid gesture from the MediaPipe landmarks, `pyautogui.click()` is triggered to "emulate a mouse click at the current cursor position".

A. Controlling Scroll by a Simple Hand Gesture (via MediaPipe Hands)

To implement scrolling, we used a "pinch-to-scroll" gesture. This gesture is detected using the MediaPipe Hands model, which tracks all the landmarks of your hand in real-time. We continuously calculate the 2D distance between the tip of the thumb (Landmark 4) and the tip of the index finger (Landmark 8):

$$d_{\text{pinch}} = \sqrt{(x_4 - x_8)^2 + (y_4 - y_8)^2} \quad (15)$$

The scroll function "turns on" if this distance becomes smaller than a pinch threshold:

$$d_{\text{pinch}} < 0.05 \quad (16)$$

Once this pinch is active, we track the vertical (Y-axis) movement of the index finger; moving the pinched hand up triggers an upward scroll, and moving it down triggers a downward scroll, both of which are executed by the `pyautogui` library.

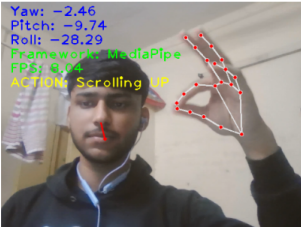


Fig. 5. Scrolling up gesture.

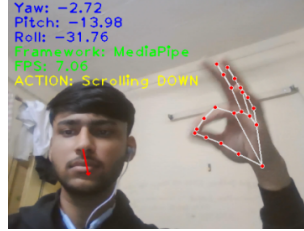


Fig. 6. Scrolling down gesture.

B. Fist Recognition (via MediaPipe Hands)

To address the limitations of eye blink detection (which we found unreliable with our low-resolution webcam), we implemented fist recognition as an alternative method for simulating mouse clicks. This approach proved to be "more reliable and easier to control in real-time".

We use the MediaPipe Hands framework, which employs the following logic:

Fist Detection (for Clicking): The code checks the y -coordinates of the landmarks. It confirms that all four fingertips (e.g., `INDEX_FINGER_TIP`) are positioned below their corresponding middle knuckles (e.g., `INDEX_FINGER_PIP`):

If this condition is true for all four fingers, the code registers a "fist."

This MediaPipe-based fist recognition offered "better stability and fewer false positives" than blink detection, all while being extremely computationally efficient.

By integrating our MediaPipe Hands detection with the `pyautogui` library, we achieved full, hands-free cursor control. Specifically, fist recognition served as a reliable



Fig. 7. Fist Gesture .

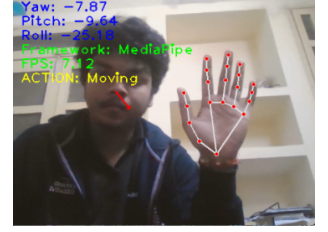


Fig. 8. Moving Gesture .

trigger for clicks, while pinch gestures provided an intuitive signal for scrolling.

VII. CHALLENGES REQUIRED

During the development of this project, we faced several significant challenges that required practical solutions.

- 1) **Initial Concept Failure (Eye-Gaze):** Our initial approach, based on the reference PDF, was to use direct eye-gaze and blink detection. This failed immediately due to the low resolution of our standard webcam. The video feed was not stable enough to track the small, rapid movements of the pupil, resulting in an unusable, "jittery" cursor. This forced us to pivot to the more robust method of head-pose estimation.

- 2) **Framework Performance (Dlib vs. MediaPipe):** Our first major challenge was selecting a viable framework. Our comparative test revealed that the Dlib implementation was unsuitable for two reasons:

- **Performance:** Dlib was computationally heavy, achieving only ~ 59 FPS, less than half of MediaPipe's ~ 123 FPS.
- **Robustness:** The `Dlib.get_frontal_face_detector()` would fail the instant the head was turned, causing all tracking to stop. MediaPipe's model, by contrast, was "capable of functioning under varying conditions, including different head poses".

- 3) **Integration Performance Cost (The "FPS Drop"):** After selecting MediaPipe as the superior framework, we faced our next major challenge. While the head-pose-only script ran at a high ~ 120 FPS, integrating the MediaPipe Hands model (in `gaze_controller.py`) caused a significant performance drop to as low as 5-15 FPS. This is because we are now running two computationally expensive deep-learning models (`face_mesh.process()` and `hands.process()`) on every single frame, which saturates the system's processing capabilities. This represents a critical trade-off between features (gestures) and real-time performance.

VIII. CONCLUSION

The comparative analysis conducted in this study demonstrates a clear and definitive result: MediaPipe is the superior framework for our real-time, head-pose-based HCI system.

Its modern, lightweight architecture places a significantly lower load on the hardware, as evidenced by its FPS being over 5 times that of Dlib's. This high efficiency, combined with its robust pose tracking, makes it the clear choice. Therefore, for all future advancements—such as integrating our YOLOv8 fist detection for clicks and MediaPipe Hands for scrolling—we will be building exclusively on the MediaPipe foundation.

Despite the limitations, our final application uses stable head-pose estimation for navigation, combined with reliable fist-detection for clicks and pinch-gestures for scrolling.

This project's success was achieved by systematically identifying and solving several key challenges. We began by pivoting from unreliable eye-tracking to stable head-pose, and then proved that MediaPipe was a superior framework to Dlib in both performance and robustness. The final system is a stable, low-latency controller that demonstrates the feasibility of a truly accessible, vision-based interface.

Future work could build on this stable foundation to include dynamic screen mapping, which would allow the controller to learn and adapt to a user's specific range of head motion.

REFERENCES

- [1] Dlib Face Landmarks Detector. Shape Predictor 68 Face Landmarks. Available: <https://www.kaggle.com/datasets/sajikim/shape-predictor-68-face-landmarks>
- [2] J. Hou, Face Yaw Roll Pitch from Pose Estimation using OpenCV. GitHub Repository. Available: https://github.com/jerryhouuu/Face-Yaw-Roll-Pitch-from-Pose-Estimation-using-OpenCV/blob/master/pose_estimation.py
- [3] Landmark Detection. Wikipedia. Available: https://en.wikipedia.org/wiki/Landmark_detection
- [4] Head Pose Estimation Using OpenCV and Dlib. Learn OpenCV. Available: <https://learnopencv.com/head-pose-estimation-using-opencv-and-dlib/>
- [5] R. Neupane, Facial Landmark Detection. Medium, 2023. Available: <https://medium.com/@RiwajNeupane/facial-landmark-detection-a6b3e29eac5b>
- [6] C. Lugaresi, J. Tang, H. Nash, M. McGuire, N. Kalantri, F. Fei, and J. Kosslyn, "MediaPipe: A Framework for Building Perception Pipelines," arXiv preprint arXiv:1906.08172, 2019. Available: <https://arxiv.org/abs/1906.08172>

APPENDIX

This appendix contains the complete Python implementation code for the head pose estimation system. Two main scripts are presented: the first enables comparison between Dlib and MediaPipe frameworks, while the second provides the final integrated system combining MediaPipe FaceMesh with MediaPipe Hands for gesture-based control.

A. Head Pose Comparison Script: Dlib vs MediaPipe

The following script implements both Dlib and MediaPipe frameworks for head pose estimation. This code demonstrates the comparative analysis discussed in Section 5 of the paper.

```
1 import cv2
2 import numpy as np
3 import dlib
4 import mediapipe as mp
5 import pyautogui
6 import time
7
8 # --- 0. CONFIGURATION ---
9 # Set this to True to use MediaPipe, False to use
10 # Dlib
11 USE_MEDIAPIPE = False
12
13 # --- 1. DLIB SETUP ---
14 print("Loading Dlib model...")
15 try:
16     dlib_detector = dlib.get_frontal_face_detector()
17     dlib_predictor = dlib.shape_predictor("shape_predictor_68_face_landmarks.dat")
18 except RuntimeError:
19     print("ERROR: Could not find 'shape_predictor_68_face_landmarks.dat'")
20     print("Download it from: http://dlib.net/files/shape_predictor_68_face_landmarks.dat.bz2")
21     print("And unzip it in the same folder as this script.")
22     exit()
23
24 # Dlib 68-point model indices
25 DLIB_LANDMARK_IDS = [
26     30, # Nose Tip
27     36, # Left Eye Outer Corner
28     45, # Right Eye Outer Corner
29     48, # Left Mouth Corner
30     54, # Right Mouth Corner
31     8, # Chin
32 ]
33
34 # 3D model points for Dlib
35 DLIB_MODEL_POINTS = np.array([
36     (0.0, 0.0, 0.0), # Nose Tip
37     (-22.0, -17.0, -26.0), # Left Eye Outer
38     (22.0, -17.0, -26.0), # Right Eye Outer
39     (-15.0, 27.0, -20.0), # Left Mouth Corner
40     (15.0, 27.0, -20.0), # Right Mouth Corner
41     (0.0, 60.0, -35.0) # Chin
42 ], dtype=np.float32)
43
44 # --- 2. MEDIAPIPE SETUP ---
45 print("Loading MediaPipe model...")
46 mp_face_mesh = mp.solutions.face_mesh
47 face_mesh = mp_face_mesh.FaceMesh(
48     static_image_mode=False,
49     max_num_faces=1,
50     refine_landmarks=True)
51
52 # MediaPipe 6-point model indices
53 MEDIAPIPE_LANDMARK_IDS = [
54     1, # Nose Tip
55     133, # Left Eye Outer Corner
56     362, # Right Eye Outer Corner
57     234, # Left Ear
58     454, # Right Ear
59     152, # Chin
60 ]
61
62 # 3D model points for MediaPipe
63 MEDIAPIPE_MODEL_POINTS = np.array([
64     (0.0, 0.0, 0.0), # Nose Tip
65     (-30.0, -30.0, -30.0), # Left Eye
66     (30.0, -30.0, -30.0), # Right Eye
```

```

65     (-60.0, 60.0, -60.0), # Left Ear
66     ( 60.0, 60.0, -60.0), # Right Ear
67     ( 0.0, 100.0, -100.0) # Chin
68 ], dtype=np.float32)
69
70 # --- 3. SHARED PARAMETERS & FUNCTIONS ---
71 # Smoothing parameters
72 alpha = 0.7
73 smoothed_yaw = None
74 smoothed_pitch = None
75
76 # Camera Parameters (will be set in functions)
77 FOCAL_LENGTH = 1
78 CENTER = (0, 0)
79
80 def rotation_matrix_to_euler_angles(R):
81     """ Converts a rotation matrix to Euler angles
82         (yaw, pitch, roll) """
83     yaw = np.arctan2(R[1, 0], R[0, 0]) * (180.0 /
84         np.pi)
85     pitch = np.arctan2(-R[2, 0], np.sqrt(R[2,
86         1]**2 + R[2, 2]**2)) * (180.0 / np.pi)
87     roll = np.arctan2(R[2, 1], R[2, 2]) * (180.0 /
88         np.pi)
89     return yaw, pitch, roll
90
91 def draw_info(img, gaze, fps, framework_name):
92     """ Draws the gaze arrow and info text on the
93         frame """
94     face = gaze["face"]
95     arrow_length = img.shape[1] / 3
96     dx = -arrow_length * np.sin(np.radians(gaze["
97         yaw"]))
98     dy = -arrow_length * np.sin(np.radians(gaze["
99         pitch"]))
100     cv2.arrowedLine(
101         img,
102         (int(face["x"]), int(face["y"])),
103         (int(face["x"] + dx), int(face["y"] + dy))
104         ,
105         (0, 0, 255), 2, cv2.LINE_AA, tipLength=0.2
106     )
107     cv2.putText(img, f"Yaw: {gaze['yaw']:.2f}",
108         (20, 30),
109         cv2.FONT_HERSHEY_PLAIN, 2, (255,
110             0, 0), 2)
111     cv2.putText(img, f"Pitch: {gaze['pitch']:.2f}"
112         , (20, 60),
113         cv2.FONT_HERSHEY_PLAIN, 2, (255,
114             0, 0), 2)
115     cv2.putText(img, f"Roll: {gaze['roll']:.2f}",
116         (20, 90),
117         cv2.FONT_HERSHEY_PLAIN, 2, (255,
118             0, 0), 2)
119
120 # Display the framework and FPS
121 color = (0, 255, 0) if framework_name == "
122     MediaPipe" else (0, 0, 255)
123 cv2.putText(img, f"Framework: {framework_name}"
124     , (20, 120),
125     cv2.FONT_HERSHEY_PLAIN, 2, color,
126     2)
127 cv2.putText(img, f"FPS: {fps:.2f}", (20, 150),
128     cv2.FONT_HERSHEY_PLAIN, 2, color,
129     2)
130
131 # --- 4. POSE ESTIMATION FUNCTIONS ---
132 def estimate_head_pose_dlib(image):
133     """ Estimates head pose using Dlib """
134     global FOCAL_LENGTH, CENTER
135     gray = cv2.cvtColor(image, cv2.COLOR_BGR2GRAY)
136     faces = dlib_detector(gray)
137
138     if not faces:
139         return None
140
141     shape = dlib_predictor(gray, faces[0])
142     ih, iw = image.shape[:2]
143
144     image_points = np.array([
145         shape.part(i).x, shape.part(i).y]

```

```

130         for i in DLIB_LANDMARK_IDS
131     ], dtype=np.float32)
132
133     # Camera Intrinsics
134     FOCAL_LENGTH = iw
135     CENTER = (iw/2, ih/2)
136     CAMERA_MATRIX = np.array([
137         [FOCAL_LENGTH, 0, CENTER[0]],
138         [0, FOCAL_LENGTH, CENTER[1]],
139         [0, 0, 1]
140     ], dtype=np.float32)
141
142     DIST_COEFFS = np.zeros((4, 1)) # No
143         distortion
144
145     # SolvePnP
146     success, rotation_vector, translation_vector =
147         cv2.solvePnP(
148             DLIB_MODEL_POINTS, image_points,
149             CAMERA_MATRIX, DIST_COEFFS
150         )
151
152     if not success:
153         return None
154
155     rotation_matrix, _ = cv2.Rodrigues(
156         rotation_vector)
157     yaw, pitch, roll =
158         rotation_matrix_to_euler_angles(
159             rotation_matrix)
160     yaw = -yaw # Adjust sign
161
162     # Get face center
163     x_coords = [p[0] for p in image_points]
164     y_coords = [p[1] for p in image_points]
165     face_x = np.mean(x_coords)
166     face_y = np.mean(y_coords)
167
168     return {"yaw": yaw, "pitch": pitch, "roll":
169         roll,
170         "face": {"x": face_x, "y": face_y}}
171
172 def estimate_head_pose_mediapipe(image):
173     """ Estimates head pose using MediaPipe """
174     global FOCAL_LENGTH, CENTER
175     img_rgb = cv2.cvtColor(image, cv2.
176         COLOR_BGR2RGB)
177     results = face_mesh.process(img_rgb)
178
179     if not results.multi_face_landmarks:
180         return None
181
182     face_landmarks = results.multi_face_landmarks
183     [0]
184     ih, iw = image.shape[:2]
185
186     # Get 2D image points from MediaPipe
187     image_points = np.array([
188         [face_landmarks.landmark[i].x * iw,
189         face_landmarks.landmark[i].y * ih]
190         for i in MEDIAPIPE_LANDMARK_IDS
191     ], dtype=np.float32)
192
193     # Camera Intrinsics
194     FOCAL_LENGTH = iw
195     CENTER = (iw/2, ih/2)
196     CAMERA_MATRIX = np.array([
197         [FOCAL_LENGTH, 0, CENTER[0]],
198         [0, FOCAL_LENGTH, CENTER[1]],
199         [0, 0, 1]
200     ], dtype=np.float32)
201
202     DIST_COEFFS = np.zeros((4, 1)) # No
203         distortion
204
205     # SolvePnP
206     success, rotation_vector, translation_vector =
207         cv2.solvePnP(
208             MEDIAPIPE_MODEL_POINTS, image_points,
209             CAMERA_MATRIX, DIST_COEFFS
210         )
211
212     if not success:

```



```

201         return None
202
203     rotation_matrix, _ = cv2.Rodrigues(
204         rotation_vector)
205     yaw, pitch, roll =
206         rotation_matrix_to_euler_angles(
207             rotation_matrix)
208     yaw = -yaw # Adjust sign
209
210     # Get face center (use nose tip)
211     face_x = image_points[0][0]
212     face_y = image_points[0][1]
213
214     return {"yaw": yaw, "pitch": pitch, "roll":
215            roll,
216            "face": {"x": face_x, "y": face_y}}
217
218 # --- 5. MAIN WEBCAM LOOP ---
219 print("Running Head Pose Comparison...")
220 print(f"*** USING {'MEDIAPIPE' if USE_MEDIAPIPE
221       else 'DLIB'} FOR POSE ESTIMATION ***")
222
223 cap = cv2.VideoCapture(0)
224 cv2.namedWindow("Head Pose Comparison")
225
226 while cap.isOpened():
227     ret, frame = cap.read()
228     if not ret:
229         break
230
231     # --- Gaze Estimation (SWITCHABLE) ---
232     framework_name = ""
233     gaze_data = None
234
235     if USE_MEDIAPIPE:
236         start_time = time.time()
237         gaze_data = estimate_head_pose_mediarpipe(
238             frame)
239         processing_time = time.time() - start_time
240         framework_name = "MediaPipe"
241     else:
242         start_time = time.time()
243         gaze_data = estimate_head_pose_dlib(frame)
244         processing_time = time.time() - start_time
245         framework_name = "Dlib"
246
247     # Calculate FPS
248     if processing_time > 0:
249         fps = 1 / processing_time
250     else:
251         fps = float('inf')
252
253     if gaze_data:
254         # Smoothing
255         if smoothed_yaw is None:
256             smoothed_yaw = gaze_data["yaw"]
257             smoothed_pitch = gaze_data["pitch"]
258         else:
259             smoothed_yaw = alpha * gaze_data["yaw"]
260             + (1 - alpha) * smoothed_yaw
261             smoothed_pitch = alpha * gaze_data["
262             pitch"] + (1 - alpha) *
263             smoothed_pitch
264
265     gaze_data["yaw"] = smoothed_yaw
266     gaze_data["pitch"] = smoothed_pitch
267
268     # Draw info
269     draw_info(frame, gaze_data, fps,
270              framework_name)
271
272     # --- PyAutoGUI Mouse Control ---
273     screen_width, screen_height = pyautogui.
274         size()
275
276     # This mapping can be tuned
277     target_x = screen_width * (gaze_data["yaw"]
278                                + 40) / 80
279     target_y = screen_height * (gaze_data["
280     pitch"] + 30) / 60
281
282     target_x = np.clip(target_x, 0,
283                        screen_width - 1)

```

```

270         target_y = np.clip(target_y, 0,
271                             screen_height - 1)
272
273         pyautogui.moveTo(target_x, target_y)
274
275         # Display the result
276         cv2.imshow("Head Pose Comparison", frame)
277
278         if cv2.waitKey(1) & 0xFF == ord("q"):
279             break
280
281     cap.release()
282     cv2.destroyAllWindows()

```

Listing 1. Head Pose Comparison: Dlib vs MediaPipe Framework

B. Final Integrated System: MediaPipe with Gesture Control

This script represents the final, optimized implementation built entirely on the MediaPipe framework.

```

1  import cv2
2  import numpy as np
3  import mediapipe as mp
4  import pyautogui
5  import time
6
7  # --- MEDIAPIPE SETUP ---
8  print("Loading MediaPipe models...")
9  mp_face_mesh = mp.solutions.face_mesh
10 mp_hands = mp.solutions.hands
11 mp_drawing = mp.solutions.drawing_utils
12
13 face_mesh = mp_face_mesh.FaceMesh(
14     static_image_mode=False,
15     max_num_faces=1,
16     refine_landmarks=True,
17     min_detection_confidence=0.5,
18     min_tracking_confidence=0.5
19 )
20
21 hands = mp_hands.Hands(
22     static_image_mode=False,
23     max_num_hands=1,
24     min_detection_confidence=0.7,
25     min_tracking_confidence=0.5
26 )
27
28 # --- CONFIGURATION ---
29 # Smoothing parameters for head pose
30 alpha = 0.7
31 smoothed_yaw = None
32 smoothed_pitch = None
33 smoothed_roll = None
34
35 # Click detection parameters
36 click_cooldown = 0.5 # Seconds between clicks
37 last_click_time = 0
38
39 # Scroll detection parameters
40 scroll_active = False
41 last_scroll_y = None
42 scroll_threshold = 0.05 # Pinch distance
43 threshold
44 scroll_sensitivity = 2
45
46 # MediaPipe 6-point model indices for head pose
47 LANDMARK_IDS = [1, 133, 362, 234, 454, 152]
48
49 # 3D model points for MediaPipe
50 MODEL_POINTS = np.array([
51     (0.0, 0.0, 0.0), # Nose Tip
52     (-30.0, -30.0, -30.0), # Left Eye
53     (30.0, -30.0, -30.0), # Right Eye
54     (-60.0, 60.0, -60.0), # Left Ear
55     (60.0, 60.0, -60.0), # Right Ear
56     (0.0, 100.0, -100.0) # Chin
57 ], dtype=np.float32)
58
59 # Camera Parameters

```

```

59 FOCAL_LENGTH = 1
60 CENTER = (0, 0)
61
62 def rotation_matrix_to_euler_angles(R):
63     """Converts a rotation matrix to Euler angles
64     (yaw, pitch, roll)"""
65     yaw = np.arctan2(R[1, 0], R[0, 0]) * (180.0 /
66         np.pi)
67     pitch = np.arctan2(-R[2, 0], np.sqrt(R[2,
68         1]**2 + R[2, 2]**2)) * (180.0 / np.pi)
69     roll = np.arctan2(R[2, 1], R[2, 2]) * (180.0 /
70         np.pi)
71     return yaw, pitch, roll
72
73 def estimate_head_pose(image):
74     """Estimates head pose using MediaPipe
75     FaceMesh"""
76     global FOCAL_LENGTH, CENTER
77     img_rgb = cv2.cvtColor(image, cv2.
78         COLOR_BGR2RGB)
79     results = face_mesh.process(img_rgb)
80
81     if not results.multi_face_landmarks:
82         return None
83
84     face_landmarks = results.multi_face_landmarks
85     [0]
86     ih, iw = image.shape[:2]
87
88     # Get 2D image points from MediaPipe
89     image_points = np.array([
90         [face_landmarks.landmark[i].x * iw,
91         face_landmarks.landmark[i].y * ih]
92         for i in LANDMARK_IDS
93     ], dtype=np.float32)
94
95     # Camera Intrinsics
96     FOCAL_LENGTH = iw
97     CENTER = (iw/2, ih/2)
98     CAMERA_MATRIX = np.array([
99         [FOCAL_LENGTH, 0, CENTER[0]],
100         [0, FOCAL_LENGTH, CENTER[1]],
101         [0, 0, 1]
102     ], dtype=np.float32)
103
104     DIST_COEFFS = np.zeros((4, 1)) # No
105     distortion
106
107     # SolvePnP
108     success, rotation_vector, _ = cv2.solvePnP(
109         MODEL_POINTS, image_points, CAMERA_MATRIX,
110         DIST_COEFFS
111     )
112
113     if not success:
114         return None
115
116     rotation_matrix, _ = cv2.Rodrigues(
117         rotation_vector)
118     yaw, pitch, roll =
119         rotation_matrix_to_euler_angles(
120         rotation_matrix)
121     yaw = -yaw # Adjust sign
122
123     # Get face center (use nose tip)
124     face_x = image_points[0][0]
125     face_y = image_points[0][1]
126
127     return {
128         "yaw": yaw,
129         "pitch": pitch,
130         "roll": roll,
131         "face": {"x": face_x, "y": face_y}
132     }
133
134 def detect_fist(hand_landmarks):
135     """Detects if hand is making a fist gesture"""
136     # Check if all fingertips are below their PIP
137     joints
138     fingers_down = []
139
140     # Index finger
141     fingers_down.append(

```

```

142         hand_landmarks.landmark[mp_hands.
143             HandLandmark.INDEX_FINGER_TIP].y >
144         hand_landmarks.landmark[mp_hands.
145             HandLandmark.INDEX_FINGER_PIP].y
146     )
147
148     # Middle finger
149     fingers_down.append(
150         hand_landmarks.landmark[mp_hands.
151             HandLandmark.MIDDLE_FINGER_TIP].y >
152         hand_landmarks.landmark[mp_hands.
153             HandLandmark.MIDDLE_FINGER_PIP].y
154     )
155
156     # Ring finger
157     fingers_down.append(
158         hand_landmarks.landmark[mp_hands.
159             HandLandmark.RING_FINGER_TIP].y >
160         hand_landmarks.landmark[mp_hands.
161             HandLandmark.RING_FINGER_PIP].y
162     )
163
164     # Pinky finger
165     fingers_down.append(
166         hand_landmarks.landmark[mp_hands.
167             HandLandmark.PINKY_TIP].y >
168         hand_landmarks.landmark[mp_hands.
169             HandLandmark.PINKY_PIP].y
170     )
171
172     return all(fingers_down)
173
174 def detect_pinch(hand_landmarks):
175     """Detects pinch gesture and returns distance
176     between thumb and index finger"""
177     thumb_tip = hand_landmarks.landmark[mp_hands.
178         HandLandmark.THUMB_TIP]
179     index_tip = hand_landmarks.landmark[mp_hands.
180         HandLandmark.INDEX_FINGER_TIP]
181
182     # Calculate Euclidean distance
183     distance = np.sqrt(
184         (thumb_tip.x - index_tip.x)**2 +
185         (thumb_tip.y - index_tip.y)**2
186     )
187
188     return distance
189
190 def draw_info(img, gaze, hand_gesture, fps):
191     """Draws information overlay on frame"""
192     # Draw gaze arrow
193     if gaze:
194         face = gaze["face"]
195         arrow_length = img.shape[1] / 3
196         dx = -arrow_length * np.sin(np.radians(
197             gaze["yaw"]))
198         dy = -arrow_length * np.sin(np.radians(
199             gaze["pitch"]))
200
201         cv2.arrowedLine(
202             img,
203             (int(face["x"]), int(face["y"])),
204             (int(face["x"] + dx), int(face["y"] +
205                 dy)),
206             (0, 0, 255), 2, cv2.LINE_AA, tipLength
207                 =0.2
208         )
209
210     # Display angles
211     cv2.putText(img, f"Yaw: {gaze['yaw']:.1f}"
212         , (20, 30),
213         cv2.FONT_HERSHEY_SIMPLEX, 0.7,
214         (255, 0, 0), 2)
215
216     cv2.putText(img, f"Pitch: {gaze['pitch']:.1f}"
217         , (20, 60),
218         cv2.FONT_HERSHEY_SIMPLEX, 0.7,
219         (255, 0, 0), 2)
220
221     # Display gesture status
222     if hand_gesture:
223         cv2.putText(img, f"Gesture: {hand_gesture}"
224             , (20, 90),
225             cv2.FONT_HERSHEY_SIMPLEX, 0.7,

```

```

(0, 255, 0), 2)
192
193 # Display FPS
194 cv2.putText(img, f"FPS: {fps:.1f}", (20, 120),
195             cv2.FONT_HERSHEY_SIMPLEX, 0.7, (0,
196             255, 255), 2)
197
198 # --- MAIN LOOP ---
199 print("Starting integrated gaze and gesture
200       control system...")
201 cap = cv2.VideoCapture(0)
202 cv2.namedWindow("Gaze Controller with Gestures")
203
204 while cap.isOpened():
205     start_time = time.time()
206     ret, frame = cap.read()
207     if not ret:
208         break
209
210     frame = cv2.flip(frame, 1) # Mirror for
211     intuitive control
212     img_rgb = cv2.cvtColor(frame, cv2.
213     COLOR_BGR2RGB)
214
215     # --- HEAD POSE ESTIMATION ---
216     gaze_data = estimate_head_pose(frame)
217
218     if gaze_data:
219         # Apply smoothing
220         if smoothed_yaw is None:
221             smoothed_yaw = gaze_data["yaw"]
222             smoothed_pitch = gaze_data["pitch"]
223         else:
224             smoothed_yaw = alpha * gaze_data["yaw"]
225             + (1 - alpha) * smoothed_yaw
226             smoothed_pitch = alpha * gaze_data["
227             pitch"] + (1 - alpha) *
228             smoothed_pitch
229
230     gaze_data["yaw"] = smoothed_yaw
231     gaze_data["pitch"] = smoothed_pitch
232
233     # --- MOUSE CURSOR CONTROL ---
234     screen_width, screen_height = pyautogui.
235     size()
236
237     # Map head pose to screen coordinates
238     target_x = screen_width * (gaze_data["yaw"
239     ] + 40) / 80
240     target_y = screen_height * (gaze_data["
241     pitch"] + 30) / 60
242
243     target_x = np.clip(target_x, 0,
244     screen_width - 1)
245     target_y = np.clip(target_y, 0,
246     screen_height - 1)
247
248     pyautogui.moveTo(target_x, target_y,
249     duration=0)
250
251     # --- HAND GESTURE DETECTION ---
252     hand_results = hands.process(img_rgb)
253     hand_gesture = None
254
255     if hand_results.multi_hand_landmarks:
256         for hand_landmarks in hand_results.
257         multi_hand_landmarks:
258             # Draw hand landmarks
259             mp_drawing.draw_landmarks(
260             frame,
261             hand_landmarks,
262             mp_hands.HAND_CONNECTIONS
263             )
264
265             # Detect fist for clicking
266             if detect_fist(hand_landmarks):
267                 current_time = time.time()
268                 if current_time - last_click_time
269                 > click_cooldown:
270                     pyautogui.click()
271                     last_click_time = current_time
272                     hand_gesture = "CLICK"
273
274

```

```

259 # Detect pinch for scrolling
260 pinch_distance = detect_pinch(
261     hand_landmarks)
262
263 if pinch_distance < scroll_threshold:
264     hand_gesture = "SCROLL"
265     index_tip = hand_landmarks.
266     landmark[mp_hands.
267     HandLandmark.INDEX_FINGER_TIP
268     ]
269
270     if last_scroll_y is not None:
271         scroll_delta = (last_scroll_y
272         - index_tip.y) *
273         scroll_sensitivity * 100
274         if abs(scroll_delta) > 10: #
275         Minimum movement
276         threshold
277         pyautogui.scroll(int(
278         scroll_delta))
279
280     last_scroll_y = index_tip.y
281 else:
282     last_scroll_y = None
283
284 # Calculate FPS
285 processing_time = time.time() - start_time
286 fps = 1 / processing_time if processing_time >
287 0 else 0
288
289 # Draw overlay
290 draw_info(frame, gaze_data, hand_gesture, fps)
291
292 # Display frame
293 cv2.imshow("Gaze Controller with Gestures",
294     frame)
295
296 if cv2.waitKey(1) & 0xFF == ord("q"):
297     break
298
299 cap.release()
300 cv2.destroyAllWindows()
301 face_mesh.close()
302 hands.close()

```

Listing 2. Final Integrated System Using MediaPipe FaceMesh and MediaPipe Hands

The Github repo for the project is <https://github.com/BlackBetty-SaysHello/Mouse-Detection-using-Head-Pose-.git>