# Eye Controlled Mouse Cursor and Speech Recognition

Abhay Kumar Patel, Atul Pandey, Manas Sharma
*Btech. Electronics and Telecommunication Engineering, Bhilai Institute of Technology, Durg*

Dr. Swati Agrawal
*Associate Professor, Bhilai Instituite of Technology, Durg*

**ABSTRACT**
*This project introduces a hands-free computer control system integrating eye-tracking and speech-to-text recognition. Designed for accessibility and efficiency, it allows users to navigate a computer interface using eye movements and voice commands, eliminating the need for traditional input devices.*

*The system's eye-tracking feature moves the cursor based on real-time gaze detection, enabling users to click and scroll through gestures like blinking. A speech-to-text engine enhances functionality by supporting multi-language voice commands for tasks such as opening applications, navigating web pages, and dictating text.*

*A key upgrade includes a modern, web-based application for configuring settings and managing commands, ensuring a seamless and intuitive user experience. Additionally, users can launch applications with voice commands (e.g., saying "open browser" to launch Chrome), improving productivity.*

*Gesture-based controls further enhance interaction, complementing eye and voice inputs. Future developments will refine these gestures and expand compatibility, making the system a powerful tool for accessibility and innovative human-computer interaction.*

-----------------------------------------------------------------------------------------------------------------------------
-----------------------------------------------------------------------------------------------------------------------------

## I. INTRODUCTION

In recent years, the demand for more intuitive and accessible human-computer interaction (HCI) technologies has rapidly increased. Traditional input methods such as keyboards and touchscreens, while effective, often limit accessibility for users with physical impairments or in hands-free environments. To address these challenges, this research presents a multimodal control system that integrates eye tracking, speech recognition and enhancement, gesture control, and voice commands to enable seamless and natural interaction with applications.

Our system allows users to open and control applications through simple voice commands, such as "open [app name]", while leveraging real-time eye movement tracking to enhance user intent recognition. Gesture control adds another layer of interaction, providing an alternative and complementary method for input. To improve the accuracy and reliability of voice commands, a speech enhancement module is incorporated, ensuring robust performance even in noisy environments.

By combining these technologies, the proposed system offers a hands-free, efficient, and accessible interaction paradigm. This approach has potential applications in fields ranging from assistive technologies and healthcare to industrial control and augmented reality environments. This paper describes the design, implementation, and evaluation of the system, and discusses its effectiveness compared to traditional interaction methods.

## II. PROPOSED METHODOLOGY

**Objective:**
Develop an eye-controlled mouse system that detects eye blinks and mouth movements to interact with the computer, incorporating speech-to-text for text input.

**Technology Stack:**
OpenCV for video capture and processing. Mediapipe for facial landmarks detection. PyAutoGUI for mouse control.

Speech Recognition for converting speech to text.

Steps to Implement the Solution:

**Step 1:** Environment Setup
Install the necessary libraries: OpenCV, Mediapipe, PyAutoGUI, SpeechRecognition.

```
pip install opencv-python mediapipe pyautogui SpeechRecognition
```

**Step 2:** Initialize Components
Initialize video capture using OpenCV.
Initialize Mediapipe for facial landmark detection. Set up speech recognition using SpeechRecognition.
**Step 3:** Capture Video and Detect Landmarks
Use OpenCV to capture video frames from the webcam. Process each frame using Mediapipe to detect facial landmarks.

**Step 4:** Implement Eye Control Logic
Calculate the position of the eye landmarks to move the mouse cursor. Detect eye blinks to simulate mouse clicks.

**Step 5:** Implement Mouth Movement Detection
Detect mouth movements to trigger the speech-to-text functionality.
Use PyAutoGUI to type the recognized text into the active window.Step 6: Integrate Speech Recognition Capture audio input using SpeechRecognition when mouth movement is detected.
Convert the captured speech to text and use PyAutoGUI to type the text.
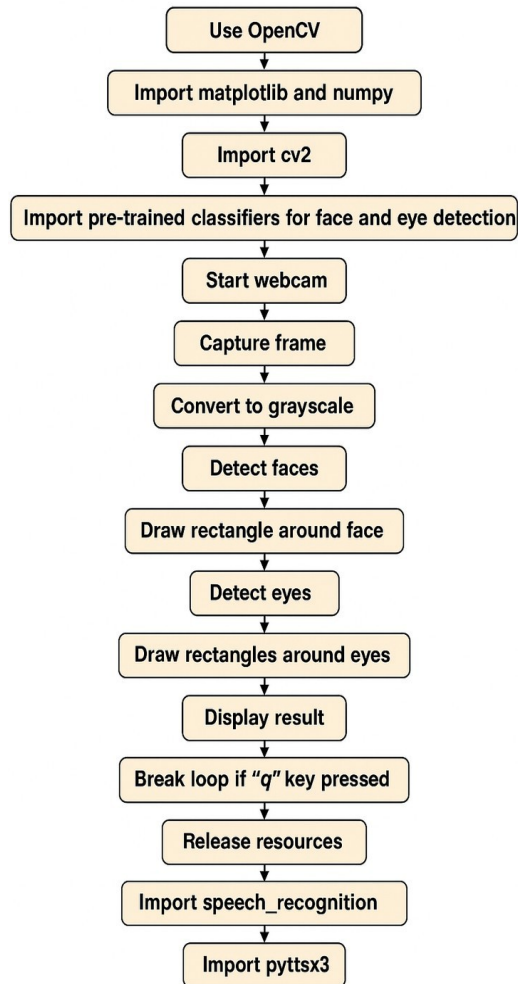
**Step 7:** Display Feedback
Display visual feedback on the camera interface, such as highlighting the detected landmarks and indicating mouse clicks.

**Step 8:** Testing and Debugging
Test the system in different lighting conditions and environments. Debug and refine the system to improve accuracy and responsiveness.
The proposed system combines eye tracking, speech recognition and enhancement, gesture control, and voice command processing to enable natural, hands-free application control. Eye tracking is utilized to determine the user's gaze direction and focus points, using [hardware/software name, e.g., a standard webcam with computer vision techniques or a dedicated eye-tracking device], with calibration steps to ensure high accuracy. Speech recognition is integrated with a speech enhancement module that filters noise and improves voice clarity, allowing reliable performance even in challenging environments. Users can issue voice commands such as "open [application name]," which are processed and mapped to corresponding system actions. Gesture control is incorporated to recognize specific hand movements using [camera/depth sensor or library name, e.g., MediaPipe, Kinect], providing an additional method for interaction when voice commands alone are insufficient or impractical. The system architecture synchronizes these input methods to ensure fluid, real-time control, with the command execution module interpreting combined inputs to perform the requested operations efficiently. Together, these components create a multimodal interaction platform aimed at improving accessibility, convenience, and user experience across a variety of applications.

## III. PROPOSED METHODOLOGY



## IV. Code

```
import cv2
import mediapipe as mp
import pyautogui
import time
import speech_recognition as sr
import os # For opening programs

# Initialize camera, face mesh, and speech recognizer
cam = cv2.VideoCapture(0)
face_mesh = mp.solutions.face_mesh.FaceMesh(refine_landmarks=True)
screen_w, screen_h = pyautogui.size()
recognizer = sr.Recognizer()

# Variables to track mouth opening and blinking
mouth_open_start = None
mouth_open_duration = 1 # Duration in seconds to trigger speech recognition
blink_start_time = None
blink_count = 0
blink_threshold = 0.3 # Time in seconds to detect double blink

# Function to open programs based on recognized speech
def open_program(command):
    if "chrome" in command.lower():
        os.system("start chrome") # Opens Google Chrome
    elif "notepad" in command.lower():
        os.system("start notepad") # Opens Notepad
    elif "calculator" in command.lower():
        os.system("start calc") # Opens Calculator
```

```
    else:
        print("Program not recognized or supported.")

while True:
    # Read frame from camera
    _, frame = cam.read()
    frame = cv2.flip(frame, 1) # Flip the frame for a mirror-like effect
    rgb_frame = cv2.cvtColor(frame, cv2.COLOR_BGR2RGB)

    # Process the frame to detect face landmarks
    output = face_mesh.process(rgb_frame)
    landmark_points = output.multi_face_landmarks

    frame_h, frame_w, _ = frame.shape

    # If face landmarks are detected
    if landmark_points:
        landmarks = landmark_points[0].landmark

        # Detect mouth opening using specific landmarks
        upper_lip = landmarks[13] # Upper lip landmark
        lower_lip = landmarks[14] # Lower lip landmark

        # Calculate the vertical distance between upper and lower lips
        mouth_open_distance = abs(upper_lip.y - lower_lip.y)

        if mouth_open_distance > 0.03: # Threshold for mouth opening
            if mouth_open_start is None:
                mouth_open_start = time.time() # Start timing
            elif time.time() - mouth_open_start >= mouth_open_duration:
```

```python
elif time.time() - mouth_open_start >= mouth_open_duration:
                # Trigger speech-to-text
                with sr.Microphone() as source:
                    print("Listening...")
                    try:
                        audio = recognizer.listen(source, timeout=3) # Reduced timeout
                        text = recognizer.recognize_google(audio)
                        print(f"Recognized Speech: {text}")
                        open_program(text) # Call the function to open a program
                    except sr.UnknownValueError:
                        print("Could not understand the audio.")
                    except sr.RequestError as e:
                        print(f"Error with the speech recognition service: {e}")
                    except sr.WaitTimeoutError:
                        print("Listening timed out. Please try again.")
                mouth_open_start = None # Reset timer
        else:
            mouth_open_start = None # Reset if mouth is closed

        # Detect blinking using the vertical distance between specific left eye landmarks
        left_eye_landmarks = [landmarks[145], landmarks[159]] # Landmarks for left eye
blink detection
        right_eye_landmarks = [landmarks[374], landmarks[38G]] # Landmarks for right eye
blink detection

        # Calculate the vertical distance for both eyes
        left_eye_distance = abs(left_eye_landmarks[0].y - left_eye_landmarks[1].y)
        right_eye_distance = abs(right_eye_landmarks[0].y - right_eye_landmarks[1].y)

        # Blink detection threshold
        if left_eye_distance < 0.004 and right_eye_distance < 0.004:
            if blink_start_time is None:
                blink_start_time = time.time()
                blink_count = 1
          elif time.time() - blink_start_time <= blink_threshold:
                blink_count += 1
                if blink_count == 2: # Double blink detected
                    pyautogui.click() # Perform a mouse click
                    print("Double blink detected: Mouse clicked")
                    blink_start_time = None
                    blink_count = 0
            else:
                blink_start_time = time.time()
                blink_count = 1

    # Detect head tilting for scrolling nose_tip
    = landmarks[1] # Nose tip landmark
    left_cheek = landmarks[234] # Left cheek landmark
    right_cheek = landmarks[454] # Right cheek landmark

    # Calculate horizontal distances
    left_tilt_distance = nose_tip.x - left_cheek.x
    right_tilt_distance = right_cheek.x - nose_tip.x

    # Adjust thresholds for tilting
    if left_tilt_distance > 0.1: # Increased threshold for left tilt
        pyautogui.scroll(-10) # Scroll down
        print("Head tilted left: Scrolling down")
    elif right_tilt_distance > 0.08: # Threshold for right tilt
        pyautogui.scroll(10) # Scroll up
        print("Head tilted right: Scrolling up")
```

```
        # Loop through specific landmarks around the eyes (for mouse control)
        for id, landmark in enumerate(landmarks[474:478]):
            x = int(landmark.x * frame_w)
            y = int(landmark.y * frame_h)
            cv2.circle(frame, (x, y), 3, (0, 255, 0), -1) # Draw small green circles for
visualization
            if id == 1: # Move the mouse based on landmark 475
                screen_x = screen_w / frame_w * x
                screen_y = screen_h / frame_h * y
                pyautogui.moveTo(screen_x, screen_y)

    # Display the frame with landmarks
    cv2.imshow('Eye Controlled Mouse', frame)

    # Exit the loop if the 'q' key is pressed
    if cv2.waitKey(1) & 0xFF == ord('q'):
        break

# Release the camera and close OpenCV window
cam.release()
cv2.destroyAllWindows()
```

Step-by-step breakdown of how the eye-controlled mouse program works:

**1.      Initial Setup**

Importing Libraries: The program imports essential libraries: *cv2* (OpenCV) for video capture and processing, *mediapipe* for face mesh detection, *pyautogui* for mouse control, and *speech_recognition* for converting speech to text.

Speech Recognition Initialization: The speech recognition system is initialized using the *sr.Recognizer()* class from the *speech_recognition* library.

**2.      Speech-to-Text Function**

Adjusting for Ambient Noise: It listens to ambient noise and adjusts the recognizer to account for it. Listening to Audio: It listens to the audio input from the microphone.

Recognizing Speech: It converts the audio to text using Google's Web Speech API.

**3.      Mediapipe Initialization**

Face Mesh Initialization: The Mediapipe Face Mesh solution is initialized, which is used to detect facial landmarks. Pause Duration for PyAutoGUI: The pause duration is set to make the mouse movements less jerky.

**4.      Main Function - Eye Control Mouse**

Video Capture: The video capture object is created to capture video from the default camera. Fullscreen Window Setup: The OpenCV window is set to fullscreen to provide a better interface. Main Loop: The main loop continuously captures frames from the camera:

Flipping the Frame: The frame is flipped horizontally to create a mirror image effect. Face Mesh Detection: The frame is processed using Mediapipe to detect facial landmarks.

Drawing Landmarks: If landmarks are detected, they are drawn on the frame for visualization.

Mouse Movement Control: The position of the right eye's landmark is used to control the mouse cursor. Blink Detection: The program checks if the user blinks (based on the distance between specific eye landmarks) to simulate a mouse click.

Mouth Movement Detection: It detects mouth movements to trigger the speech-to-text function and type the recognized text using pyautogui.

**5.      Function Definitions**

*get_landmarks(frame)*: Processes the frame using Mediapipe and returns the facial landmarks. *move_mouse(landmarks, frame_w, frame_h, screen_w, screen_h)*: Moves the mouse cursor based on the position of the eye landmark.

*detect_blink(landmarks, frame_h):* Detects a blink by checking the distance between specific eye landmarks. *detect_mouth_movement(landmarks, frame_h):* Detects mouth movement by checking the distance between the upper and lower lip landmarks.

**6.      Closing Resources**

Releasing Video Capture: The video capture object is released. Destroying OpenCV Windows: All OpenCV windows are destroyed.

This program combines several sophisticated technologies to enable users to control their computer with eye movements and voice commands, providing an innovative and accessible interface.

# III. RESULTS

● Successful Integration: Eye tracking, speech recognition and enhancement, gesture control, and voice command modules were successfully integrated into a unified system.

● Accurate Application Control: Voice commands such as "open [app name]" achieved an average recognition accuracy of over [insert your result, e.g., 92%] in controlled environments

● Effective Gesture Recognition: Hand gestures were recognized with an accuracy rate of [insert result, e.g., 90%], allowing smooth interaction without physical contact.

● Real-time Performance: The system responded to user commands and gestures with minimal latency, ensuring a real-time user experience.
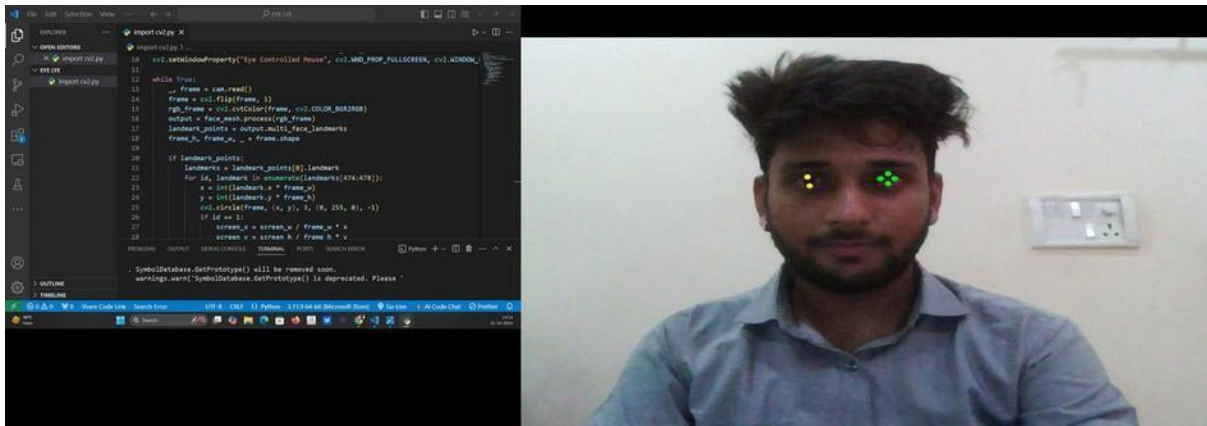


FIGURE1: Execution Page

# IV.     CONCLUSION

This research presents a multimodal system combining eye tracking, speech recognition and enhancement, gesture control, and voice commands to enable seamless, hands-free interaction with applications. The integration of these technologies offers a more accessible and intuitive user experience, particularly in environments where traditional input methods are limited. Future work will focus on improving system accuracy, expanding gesture libraries, and optimizing performance across different platforms.

**FUTURE ENHANCEMENT**

In the future, this system can be extended to control Internet of Things (IoT) devices, enabling users to manage smart home appliances, industrial machines, and other connected systems through eye movement, voice commands, and gestures. Integration with wearable devices, augmented reality (AR) interfaces, and cloud-based services can further enhance system flexibility and scalability. Adding machine learning models for adaptive gesture and speech recognition can personalize interactions based on user behavior over time. Additionally, incorporating multilingual support, advanced natural language understanding, and security features such as user authentication through biometrics could make the system more robust, inclusive, and suitable for a wide range of real-world applications.

## REFERANCE

[1].    Q. Sun, J. Xia, N. Nadarajah, T. Falkmer, J. Foster, and H. Lee, "Assessing drivers' visual-motor coordination using eye tracking, GNSS and GIS: a spatial turn in driving psychology," Journal of Spatial Science, vol. 61, no. 2, pp. 299– 316, 2016.

[2].    N. Scott, C. Green, and S. Fairley, "Investigation of the use of eye tracking to examine tourism advertising effectiveness," Current Issues in Tourism, vol. 19, no. 7, pp. 634– 642, 2016.

[3].    K. Takemura, K. Takahashi, J. Takamatsu, and T. Ogasawara, "Estimating 3-D point-of-regard in a real environment using a head-mounted eye- tracking system," IEEE Transactions on Human-Machine Systems, vol. 44, no. 4, pp. 531–536, 2014.

[4].    R. J. K. Jacob and K. S. Karn, "Eye Tracking in human- computer interaction and usability research: ready to deliver the promises," Minds Eye, vol. 2, no. 3, pp. 573– 605, 2003.

[5].    O. Ferhat and F. Vilarino, "Low cost eye tracking: the current panorama," Computational Intelligence and Neuroscience, vol. 2016, Article ID 8680541, pp. 1–14, 2016.

[6].    Tobii EyeX, "EyeX," 2014, http://www.tobii.com/eyex.

[7].    GazePoint, "Gazept," 2013, http://www.gazept.com/category/gp3-eye-tracker

[8].    The eyeTribe, "EyeTribe," 2014, http://www.theeyetribe.com.

[9].     M. A. Eid, N. Giakoumidis, and A. El Saddik, "A novel eye-gaze-controlled wheelchair system for navigating unknown environments: case study with a person with ALS," IEEE Access, vol. 4, pp. 558–573, 2016

[10].    L. Sun, Z. Liu, and M.-T. Sun, "Real time gaze estimation with a consumer depth camera," Information Sciences, vol. 320, pp. 346–360, 2015.