# IntelleX- A Multi Model Personal AI Assistant

**Deep Neural Networks**

*Submitted by*

## Navroop, Shail Sharma

**2210993719, 2210993837**

**Semester: 6th**

**BE-CSE  (Artificial  Intelligence)**

*Submitted To*

**Mr. Harsh Dalal**

Prof, Department of CSE(AI),

CUIET, Chitkara University

**CHITKARA  UNIVERSITY  INSITUTE  OF  ENGINEERING &  TECHNOLOGY**

**CHITKARA UNIVERSITY, RAJPURA**

**June, 2025**

# ACKNOWLEDGEMENTS

With immense please I, Mr/Ms Navroop (2210993819) & Shail (2210993837) presenting "IntelleX- A Multi Model Personal AI Assistant" project report as part of the curriculum of 'BE-CSE (AI)'.

I would like to express my sincere thanks to **Dr. Manu Midha and Mr Harsh Dalal**, for her & his valuable guidance and support in completing my project.

I would also like to express my gratitude towards our dean **Dr. Sushil Kumar Narang** giving me this great opportunity to do a project on NLP, MV & GEN-AI Without their support and suggestions, this project would not have been completed.

Signature………..

Name: Navroop

Roll No: 2210993819

Signature………..

Name: Shail

Roll No: 2210993837

# Table of Contents

# 1.Introduction

In the modern era, technology is evolving at an unprecedented pace, influencing every aspect of our lives, from communication and education to entertainment and productivity. One such innovation is **IntelleX**, an AI-powered virtual assistant designed to empower users through multi-modal interaction using advanced **voice commands**, **gesture recognition**, and **automated system control**. Built with tools like **Python**, **Flask**, **OpenCV**, and various AI/ML libraries, IntelleX represents a seamless bridge between human intent and machine action.

The motivation behind IntelleX lies in the growing need for more **natural, efficient, and accessible interfaces** to interact with computers. Traditional input methods often lack inclusivity and can be limiting for users with disabilities or for scenarios that demand hands-free operation. IntelleX addresses this by offering an intelligent assistant capable of executing tasks like opening applications, playing music, sending emails, browsing the internet, and more—all through voice and gesture input.

In academic and professional environments, where time efficiency and multitasking are crucial, such a system offers significant advantages. Many existing virtual assistants are either cloud-dependent or limited in scope. IntelleX is designed to be **modular, local-hosted, and privacy-conscious**, offering customizable functionalities that adapt to various user contexts. Its hands-free interface is particularly beneficial for those who work in dynamic or multitasking setups, including developers, educators, and people with mobility challenges.

This project combines core concepts from **natural language processing (NLP)**, **computer vision**, and **machine learning**, enabling it to understand spoken language, recognize gestures, and perform predefined system actions accordingly. The integration of a gesture recognition module using **Mediapipe**, and voice command processing through **SpeechRecognition** and **pyttsx3**, provides a fluid and interactive experience.

The robust architecture of IntelleX leverages **Flask** for backend control, while Python modules handle task execution, voice processing, and vision-based gesture tracking. The system is scalable and extensible, capable of integrating additional features like calendar management, task lists, system information, weather APIs, and even mobile notifications.

As the demand for intelligent personal assistants continues to grow, IntelleX showcases the practical application of AI in everyday tasks and offers a foundation for future innovations in the domain of human-computer interaction. This report presents a comprehensive overview of the methodologies, technologies, and features implemented in IntelleX, along with potential avenues for enhancement and real-world deployment.

# 2.Methodology

The IntelleX project employs a user-focused approach to create a multi-modal virtual assistant that integrates gesture recognition, voice commands, and intelligent system control. The methodology ensures seamless interaction between the user and technology, offering intuitive features designed to maximize accessibility, productivity, and hands-free convenience.

## Core Functionalities

1. **Gesture Recognition**
   - Leveraging computer vision libraries such as **OpenCV** and **Mediapipe**, IntelleX detects and interprets specific hand gestures to trigger actions like media playback, volume adjustment, brightness control, screen lock/unlock, and launching applications.
   - This functionality enables a **contactless control experience**, improving accessibility for users with physical constraints and enhancing overall system interaction.
2. **Voice Commands**
   - Using a combination of **SpeechRecognition**, **pyttsx3**, and **Flask**, the system processes spoken user instructions and responds with synthesized speech. Users can execute commands such as:
     - Opening websites and applications
     - Playing music or accessing files
     - Performing Google or Wikipedia searches
     - Sending emails or fetching weather updates
   - The voice module is tightly integrated into the backend to ensure **real-time interaction and accurate command mapping**.
3. **Object Detection** *(Upcoming Feature)*
   - The assistant is being extended to include **object detection capabilities** using deep learning frameworks (e.g., YOLOv5 or MobileNet).
   - This functionality will allow the system to identify and describe real-world objects via the camera, aiding users in tasks such as object recognition, accessibility navigation, or inventory identification.

## System Design

- **Backend Development**
  The system's core is built using **Flask**, which facilitates modular API integration, voice and gesture routing, and communication between the frontend and functional modules. Flask ensures that each command is processed efficiently and the appropriate response is delivered without delay.
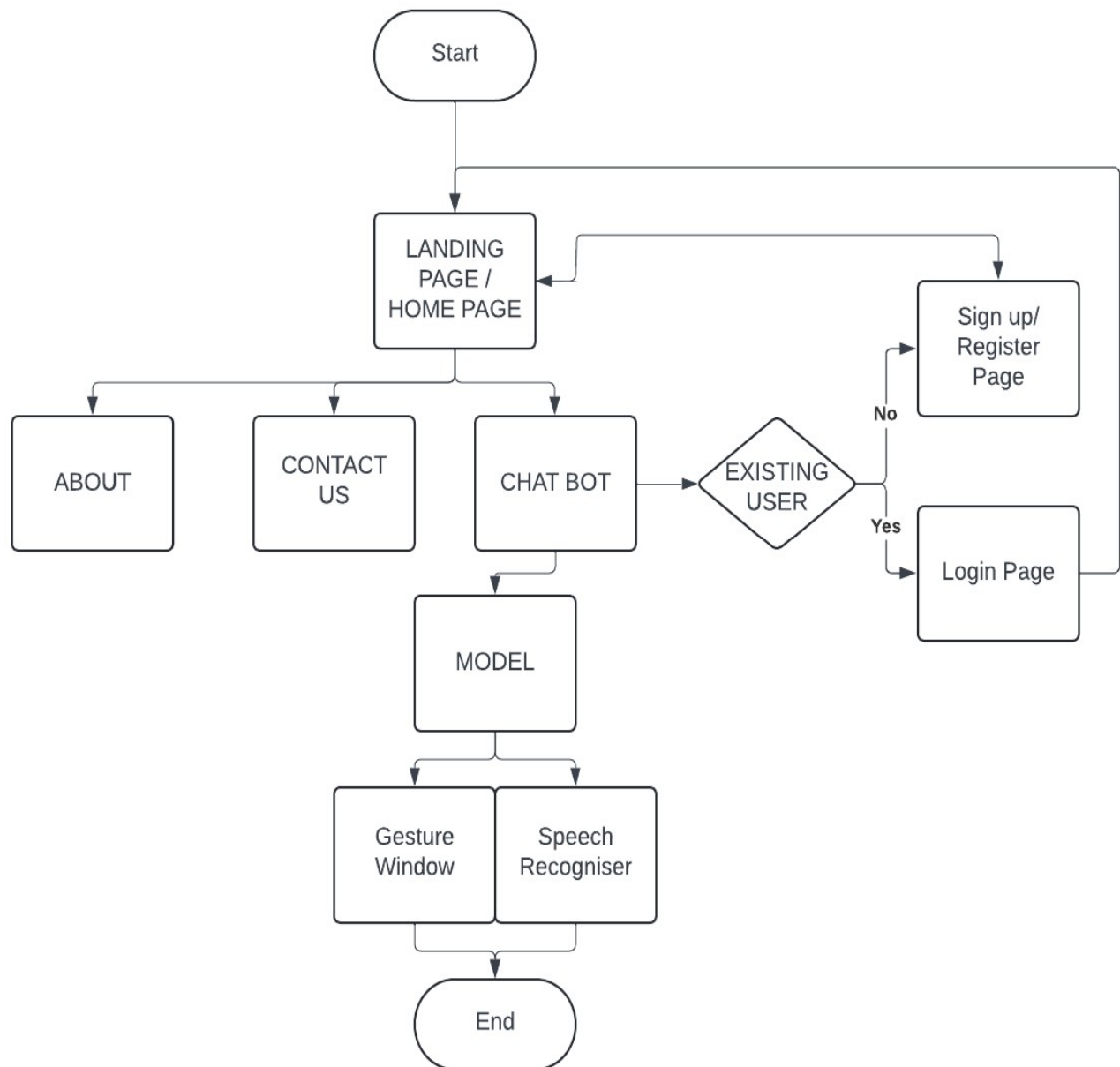- **Multi-modal Input Integration**
  IntelleX is designed to support **simultaneous visual and auditory inputs**, allowing users to interact with the assistant via hand gestures, voice, or a combination of both. This **versatile interaction model** ensures that the platform remains adaptable to different user preferences and situational needs.
- **Real-Time Feedback**
  All command responses—whether visual, verbal, or functional—are executed in **real time**, ensuring that users receive instant feedback. This creates a **fluid and interactive experience** where users can engage with their systems effortlessly and efficiently.

# 3.FlowChart

# 4.Tools and Technologies

## Frameworks and Libraries

- **Flask**
  A lightweight web framework used to build the backend of IntelleX, manage API routing, and handle communication between modules such as gesture and voice control systems.
- **OpenCV**
  An open-source computer vision library that enables real-time video processing and gesture tracking through the webcam.
- **Mediapipe**
  A Google-developed framework for building perception pipelines, used in IntelleX for accurate hand gesture and facial landmark detection.
- **Pycaw**
  A Python library that allows control over system audio, enabling users to adjust the volume through hand gestures.
- **SpeechRecognition**
  Enables speech-to-text conversion, allowing IntelleX to understand and process spoken commands.
- **pyttsx3 (Text-to-Speech)**
  A text-to-speech conversion library used to give verbal feedback to users after executing commands.
- **TensorFlow / Keras**
  Machine learning frameworks used (or to be used in future extensions) for building and deploying deep learning models such as object detection using convolutional neural networks (CNNs).
- **NumPy**
  A numerical computing library that supports image processing tasks and underpins various AI-related computations.
- **Matplotlib**
  A data visualization tool used during testing and debugging for plotting results like model accuracy, confusion matrices, and performance comparisons.

## Development Tools

- **Visual Studio Code (VS Code)**
  The primary integrated development environment (IDE) used for coding, testing, and managing the entire project.
- **Jupyter Notebook**
  Used for interactive experimentation and debugging, especially during the training and evaluation of ML models.

## Hardware Components

- **Webcam**
  Captures real-time video input for gesture recognition and object detection functionalities.
- **Microphone**
  Captures user voice input, allowing for real-time speech recognition and interaction with the assistant.
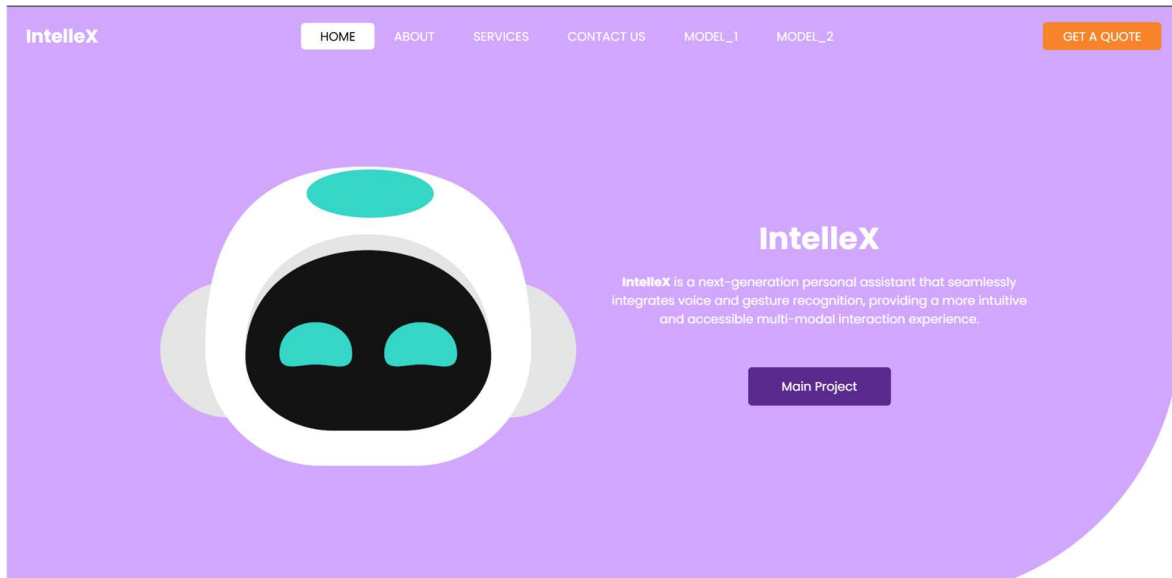
# 5.Major Findings/Outcomes/Output/Results



Fig: - 1(The image represents the homepage of "IntelleX")
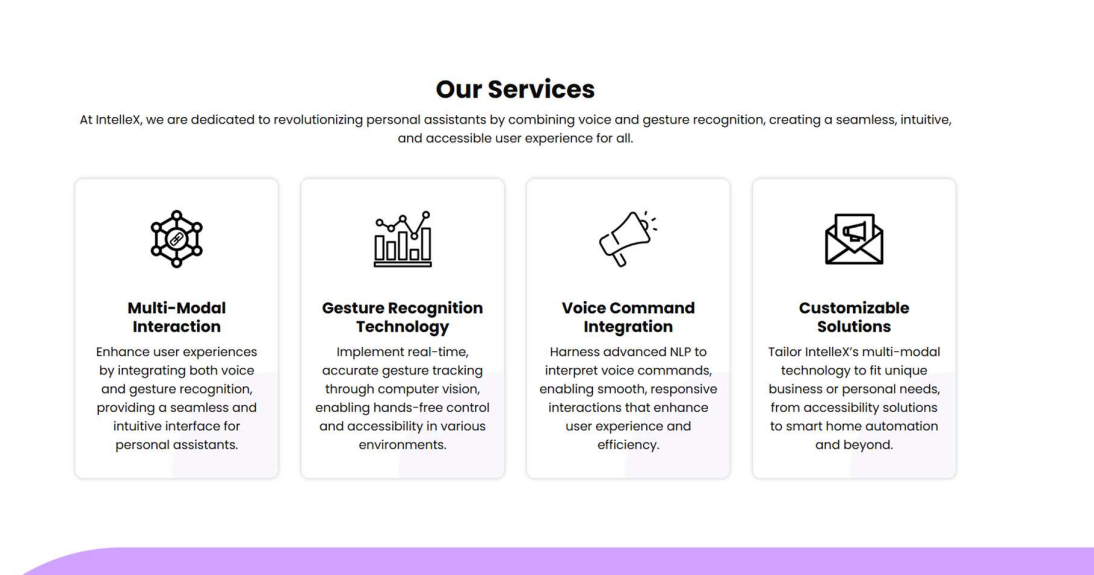


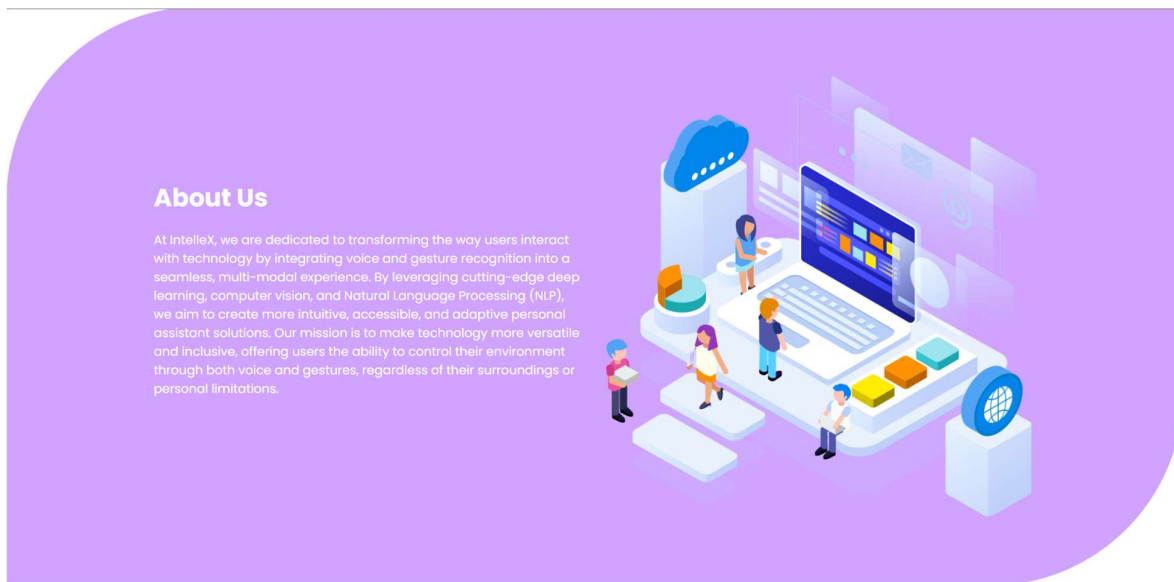Fig: - 2(The image represents the our Services of "IntelleX")

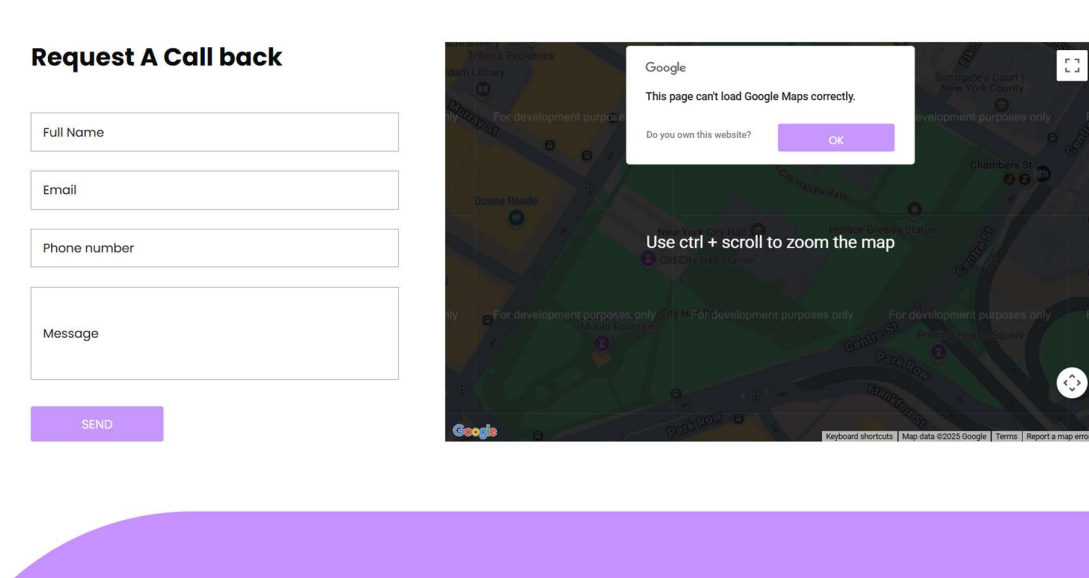Fig: - 3 (The image represents the About Us of "IntelleX")



Fig: - 4 (The image represents the footer of "IntelleX")-1
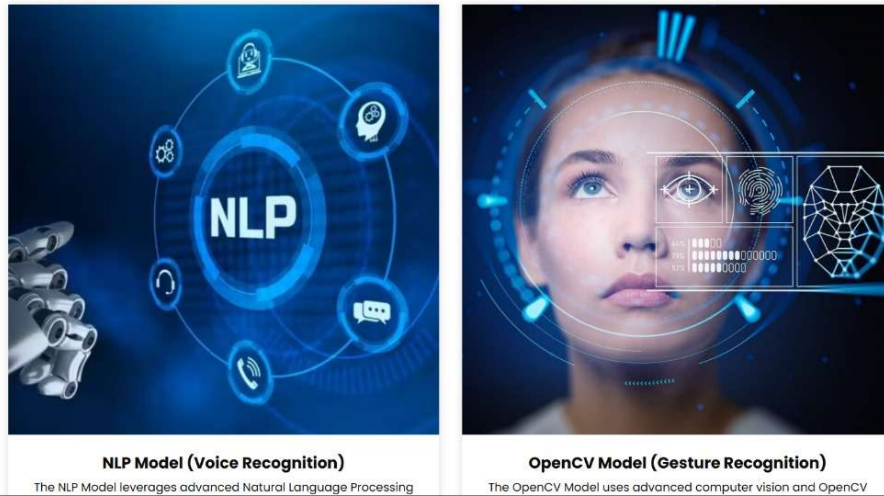
**Our Two Types Of Models**



**NLP Model (Voice Recognition)**
The NLP Model leverages advanced Natural Language Processing

**OpenCV Model (Gesture Recognition)**
The OpenCV Model uses advanced computer vision and OpenCV

Fig: - 5(The image represents the types of models of "IntelleX") -1



**NLP Model (Voice Recognition)**
The NLP Model leverages advanced Natural Language Processing (NLP) to accurately interpret voice commands and provide dynamic, efficient responses. This model is designed for users who prefer voice interactions, allowing them to control devices, access information, or perform tasks simply by speaking. It ensures a smooth and responsive user experience in a variety of environments.

**OpenCV Model (Gesture Recognition)**
The OpenCV Model uses advanced computer vision and OpenCV techniques to recognize and interpret gestures in real time with high accuracy. This model enables users to control devices and perform actions through hand movements, facial expressions, or body gestures, making it ideal for hands-free environments. It enhances accessibility and provides an intuitive, physical, and seamless way to interact with technology.

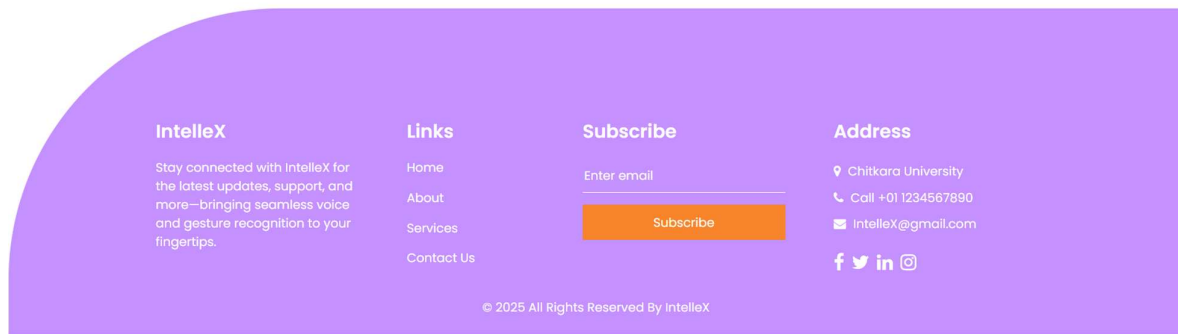Fig: - 6(The image represents the types of models of "IntelleX") -2

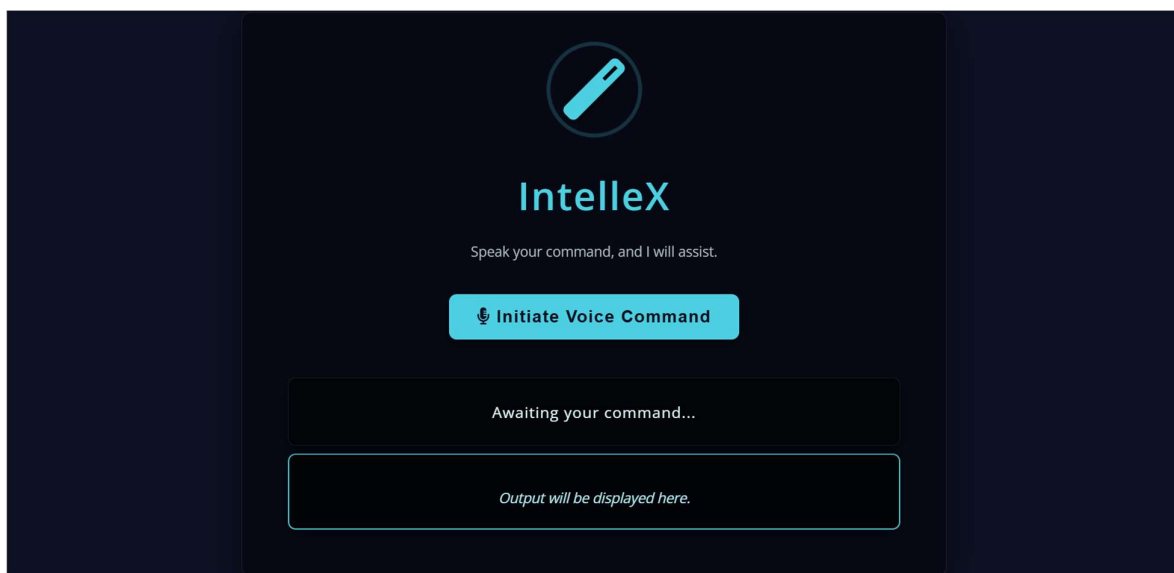Fig: - 7(The image represents the footer of "IntelleX") -2



Fig: - 8(The image represents the main Project of "IntelleX")



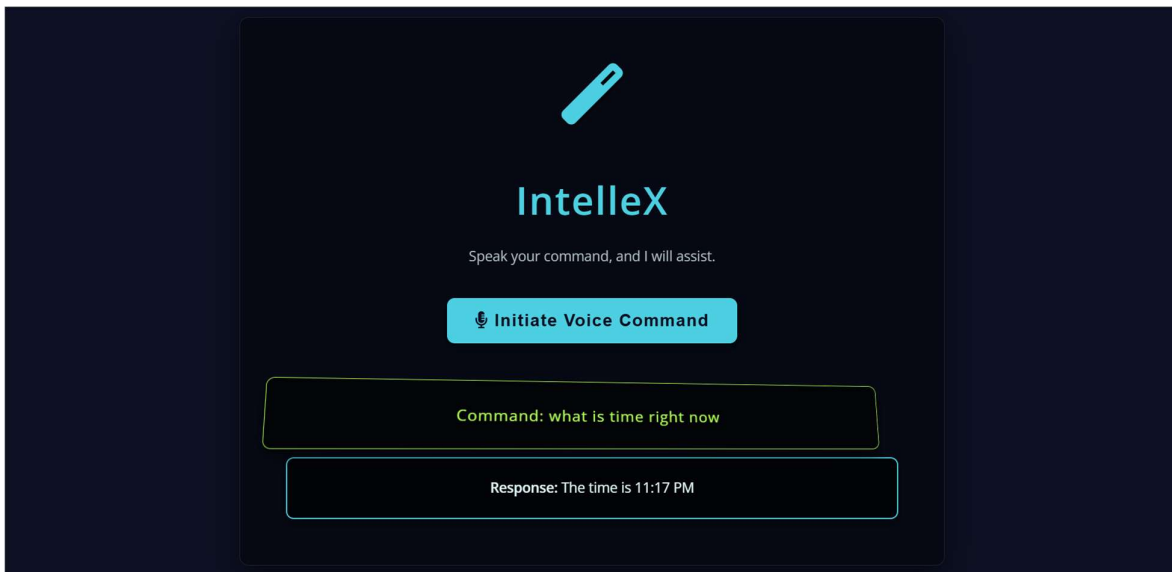Fig: - 9(The image represents the working of main Project of "IntelleX") - 1

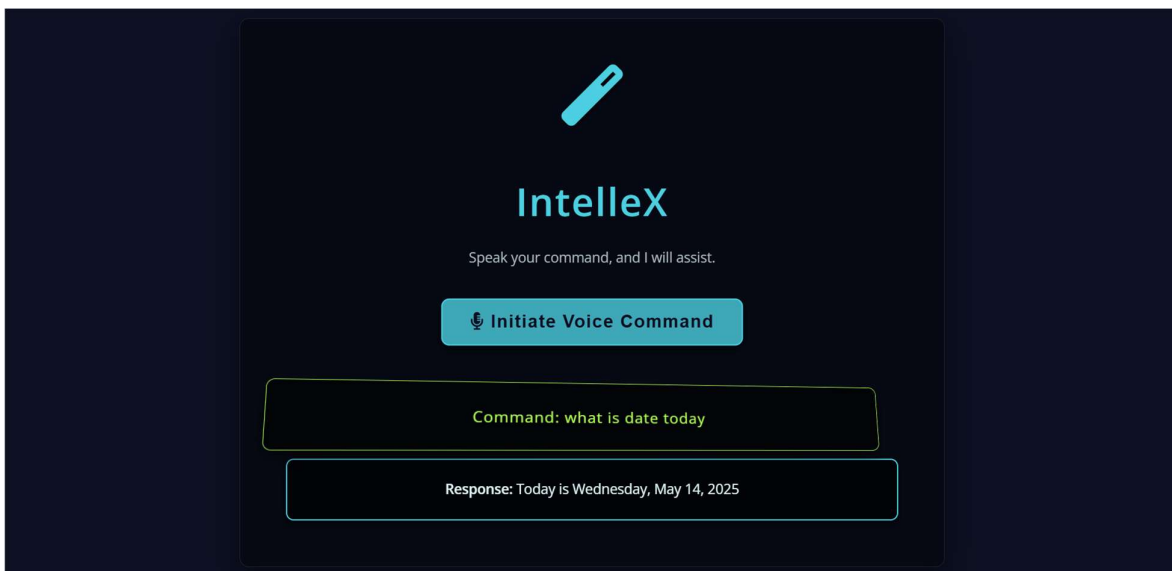Fig: - 10(The image represents the working of main Project of "IntelleX") – 2



Fig: - 11(The image represents the working of main Project of "IntelleX") – 3

Fig: - 12(The image represents the working of main Project of "IntelleX") – 4



Fig: - 13(The image represents the working of main Project of "IntelleX") – 5
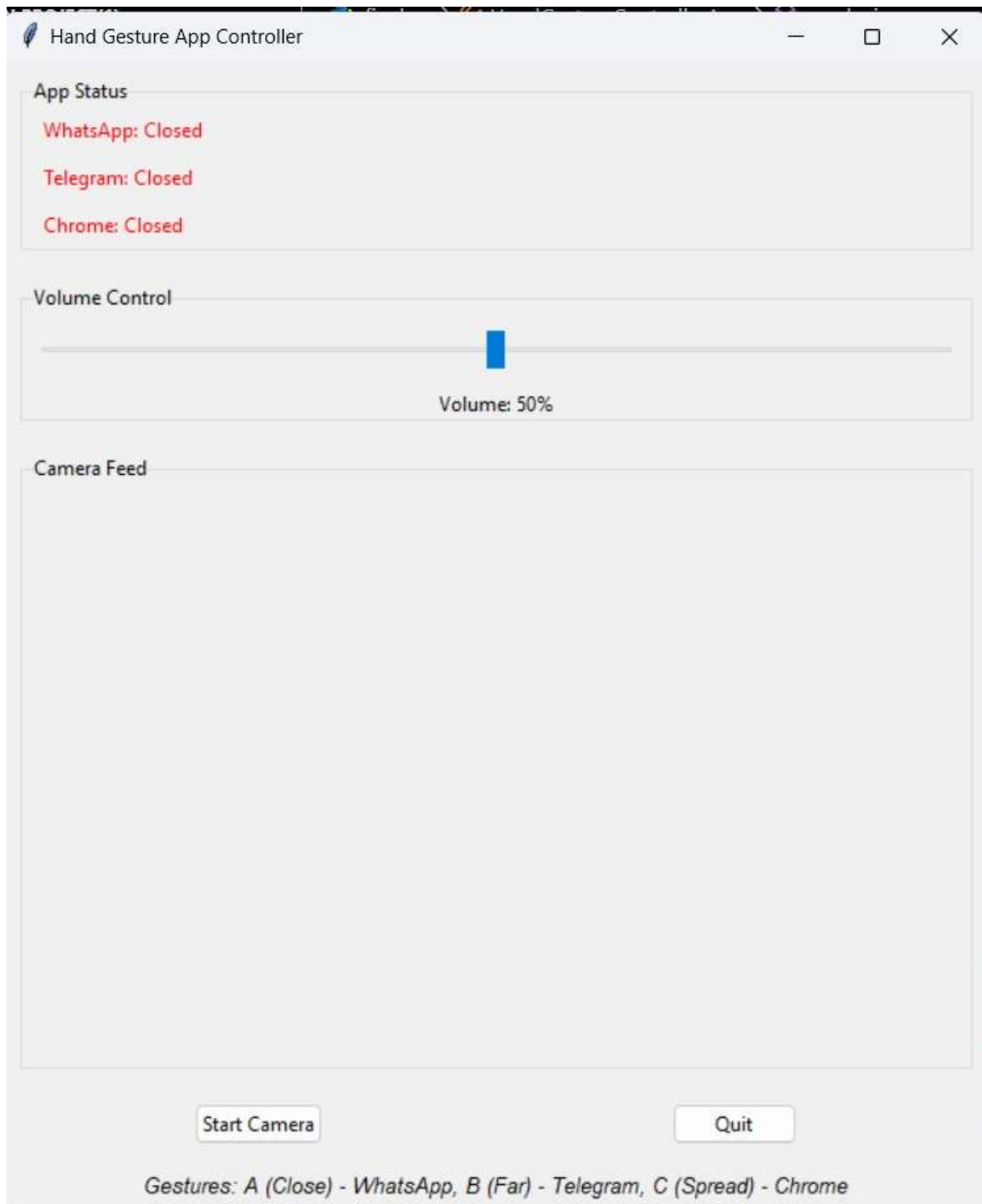
Fig: - 14(The image represents the working of main Project of "IntelleX") – 6
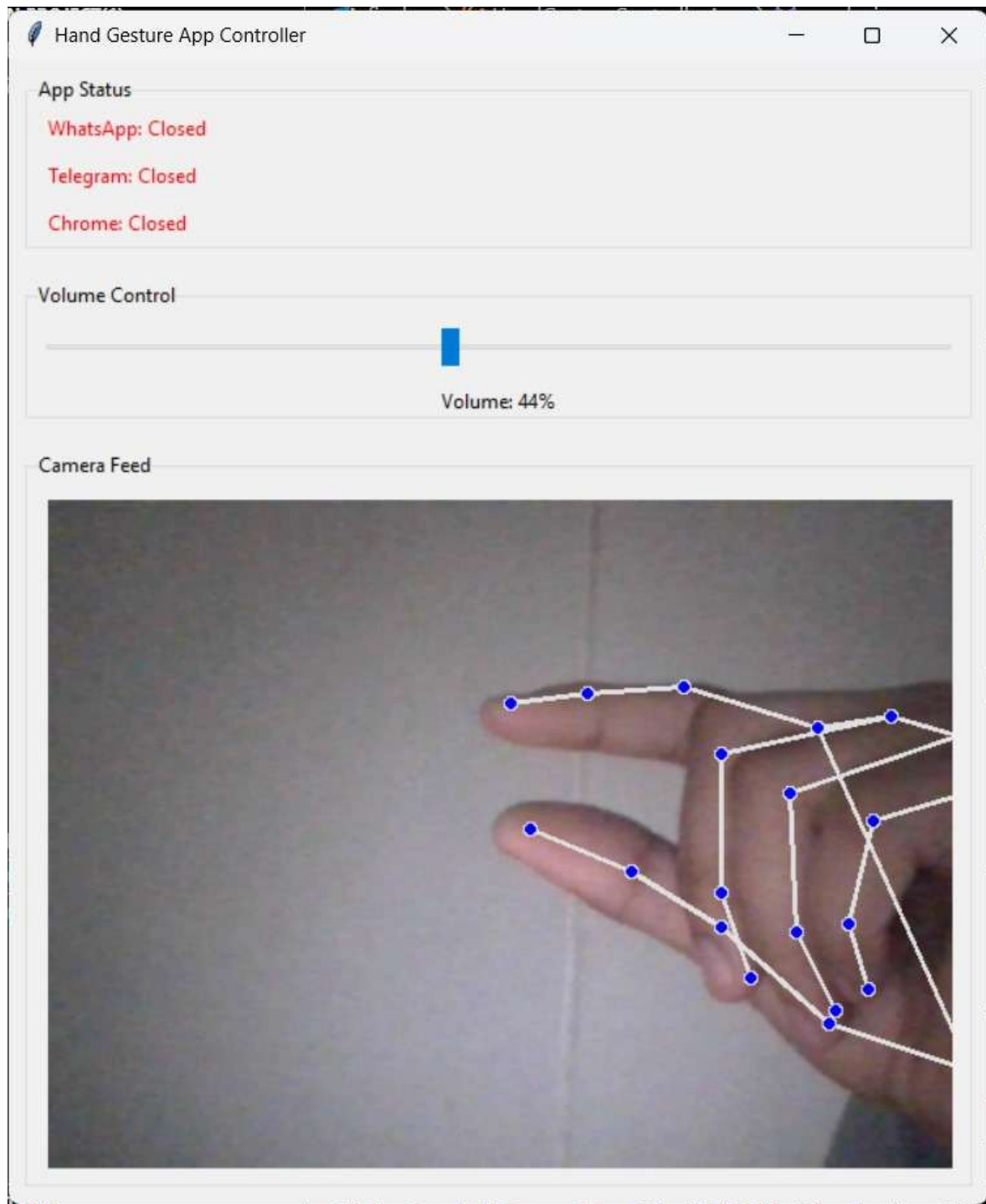
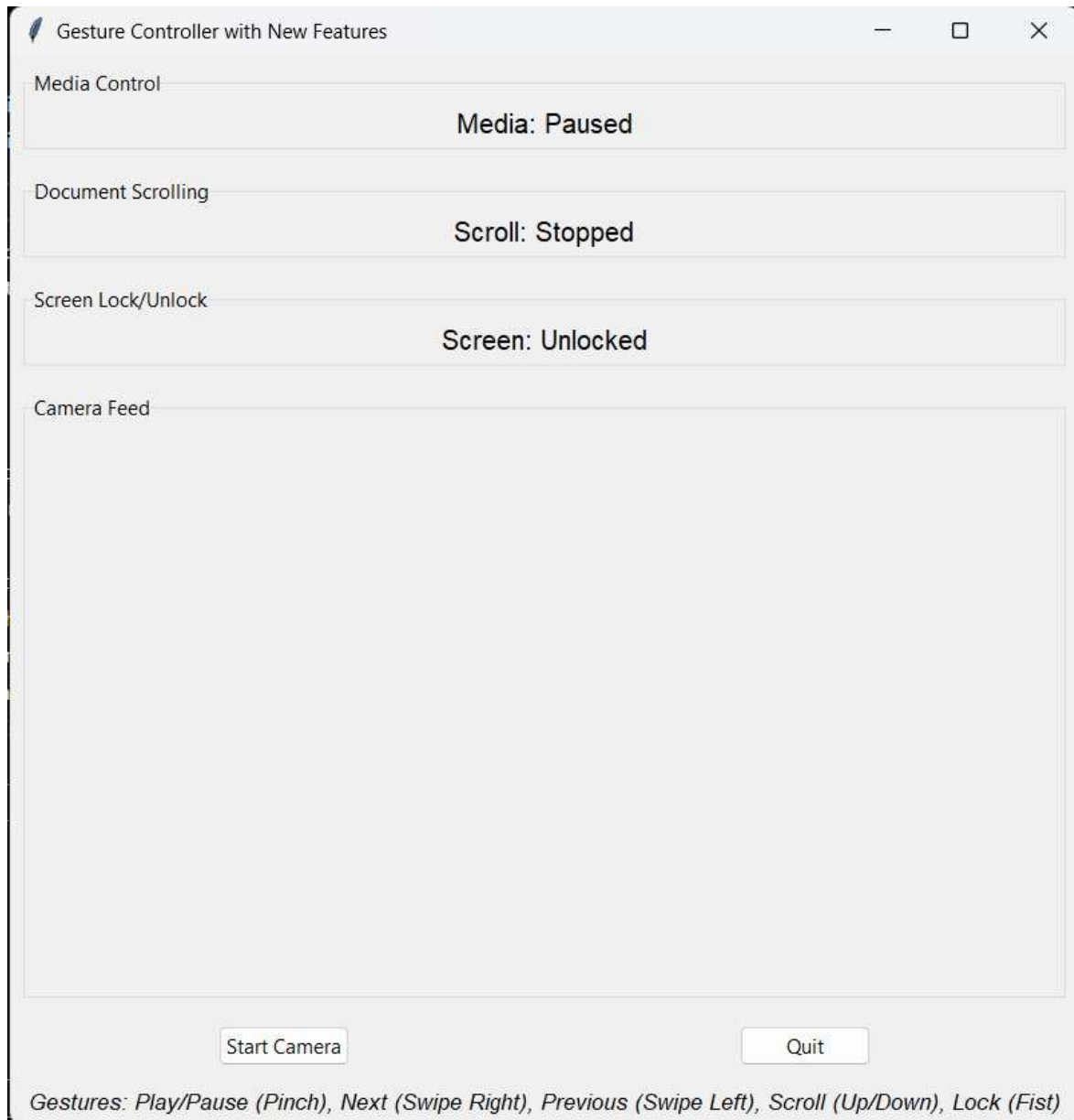Fig: - 15(The image represents the working of main Project of "IntelleX") – 7

Fig: - 16(The image represents the working of main Project of "IntelleX") – 8

```
import os
import tensorflow as tf
from tensorflow.keras.preprocessing.image import ImageDataGenerator
from tensorflow.keras.layers import Conv2D, MaxPooling2D, Flatten, Dense, Dropout
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import GlobalAveragePooling2D, Dense, Dropout
from tensorflow.keras.optimizers import Adam
from tensorflow.keras.callbacks import EarlyStopping
import matplotlib.pyplot as plt
from tensorflow.keras.preprocessing.image import load_img, img_to_array
import pandas as pd
import numpy as np
import warnings
warnings.filterwarnings("ignore")
# Disable XLA
os.environ['TF_XLA_FLAGS'] = '--tf_xla_enable_xla_devices=false'

# Set logging level to avoid unnecessary warnings
os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
tf.get_logger().setLevel('ERROR')

# Clear TensorFlow session
from tensorflow.keras import backend as K
```

Fig: - 17(The image represents the code of main Project of "IntelleX") – 1

```
# Label map
gesture_folders = [
    '01_palm', '02_l', '03_fist', '04_fist_moved', '05_thumb',
    '06_index', '07_ok', '08_palm_moved', '09_c', '10_down'
]
label_map = {gesture: idx for idx, gesture in enumerate(gesture_folders)}

# Display the images
display_images(file_paths, labels, label_map)
```
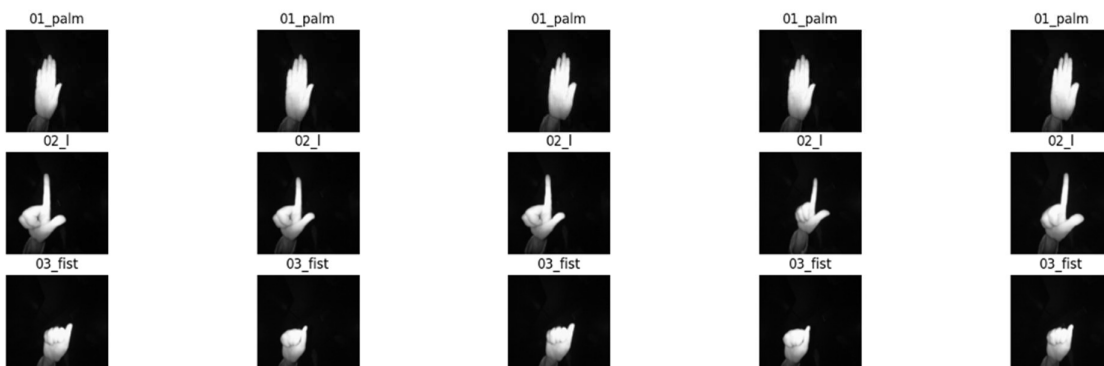


Fig: - 18(The image represents the code of main Project of "IntelleX") – 2

```
1   from flask import Flask, jsonify, render_template, session
2   import speech_recognition as sr
3   import pyttsx3
4   import webbrowser
5   import datetime
6   import os
7   import random
8   import subprocess  # To handle opening applications
9   import smtplib  # For sending emails
10  from email.mime.text import MIMEText
11  from email.mime.multipart import MIMEMultipart
12
13  app = Flask(__name__)
14  app.secret_key = 'rushi'
15
16  # Initialize the text-to-speech engine
17  engine = pyttsx3.init()
18
19  # Get all available voices
20  voices = engine.getProperty('voices')
21
22  # Set the default voice to female (session or default fallback)
23  def get_voice_id():
24      return session.get('voice_id', voices[1].id)
25
26  def speak(text, speed=130, voice_id=None):
27      if voice_id is None:
28          voice_id = get_voice_id()
29      engine.setProperty('rate', speed)
30      engine.setProperty('voice', voice_id)
31      engine.say(text)
32      engine.runAndWait()
33
34  # Recognize speech input from microphone
35  def recognize_speech():
36      recognizer = sr.Recognizer()
37      with sr.Microphone() as source:
38          print("Listening...")
```

**Fig: - 19** (The image represents the code of main Project of "IntelleX") – 3

```
1
2   import tkinter as tk
3   from tkinter import ttk
4   import cv2
5   import mediapipe as mp
6   import numpy as np
7   from AppOpener import open
8   from pycaw.pycaw import AudioUtilities, IAudioEndpointVolume
9   from ctypes import cast, POINTER
10  from comtypes import CLSCTX_ALL
11  import threading
12  import PIL.Image, PIL.ImageTk
13
14  class HandGestureControllerApp:
15      def __init__(self, root):
16          self.root = root
17          self.root.title("Hand Gesture App Controller")
18          self.root.geometry("600x700")
19
20          # Mediapipe setup
21          self.mp_hands = mp.solutions.hands
22          self.mp_drawing = mp.solutions.drawing_utils
23
24          # Apps dictionary
25          self.apps = {
26              'A': 'whatsapp',
27              'B': 'telegram',
28              'C': 'chrome'
29          }
30
31          # Volume control setup
32          devices = AudioUtilities.GetSpeakers()
33          interface = devices.Activate(IAudioEndpointVolume._iid_, CLSCTX_ALL, None)
34          self.volume = cast(interface, POINTER(IAudioEndpointVolume))
35          self.minVol, self.maxVol, _ = self.volume.GetVolumeRange()
36
37          # App state variables
38          self.last_opened_app = None
39          self.last_gesture = None
```

**Fig: - 20** (The image represents the code of main Project of "IntelleX") – 4

```python
class HandGestureControllerApp:
    def __init__(self, root):
        self.root = root
        self.root.title("Gesture Controller with Hand Tracking")
        self.root.geometry("800x800")

        # Mediapipe setup
        self.mp_hands = mp.solutions.hands.Hands(
            model_complexity=1,
            min_detection_confidence=0.7,
            min_tracking_confidence=0.7,
            max_num_hands=1
        )
        self.mp_drawing = mp.solutions.drawing_utils

        # Volume control setup
        devices = AudioUtilities.GetSpeakers()
        interface = devices.Activate(IAudioEndpointVolume._iid_, CLSCTX_ALL, None)
        self.volume = cast(interface, POINTER(IAudioEndpointVolume))
        self.minVol, self.maxVol, _ = self.volume.GetVolumeRange()

        # App state variables
        self.last_gesture = None
        self.consecutive_gesture_frames = 0

        # Camera and recognition state
        self.camera_active = False
        self.cap = None

        # Create GUI components
        self.create_widgets()

    def create_widgets(self):
        # Frame for media control
        media_frame = ttk.LabelFrame(self.root, text="Media Control")
        media_frame.pack(padx=10, pady=10, fill='x')

        self.media_status = ttk.Label(media_frame, text="Media: Paused", font=('Arial', 12))
```

Fig: - 21 (The image represents the code of main Project of "IntelleX") – 5

```python
app = Flask(__name__)
app.secret_key = 'rushi'

# Initialize the text-to-speech engine
engine = pyttsx3.init()

# Get all available voices
voices = engine.getProperty('voices')

# Set the default voice to female (session or default fallback)
def get_voice_id():
    return session.get('voice_id', voices[1].id)

def speak(text, speed=130, voice_id=None):
    if voice_id is None:
        voice_id = get_voice_id()
    engine.setProperty('rate', speed)
    engine.setProperty('voice', voice_id)
    engine.say(text)
    engine.runAndWait()

# Recognize speech input from microphone
def recognize_speech():
    recognizer = sr.Recognizer()
    with sr.Microphone() as source:
        print("Listening...")
        recognizer.adjust_for_ambient_noise(source)
        audio = recognizer.listen(source)

        try:
            query = recognizer.recognize_google(audio)
            return query.lower()
        except sr.UnknownValueError:
            speak("Sorry, I couldn't understand what you said.")
            return None
        except sr.RequestError:
```

Fig: - 22(The image represents the code of main Project of "IntelleX") – 6

# 6.Conclusion

IntelleX is a groundbreaking project that leverages advanced computer vision and machine learning techniques to deliver a multi-modal, intuitive user experience. By integrating gesture recognition, voice commands, and real-time object detection, IntelleX provides a hands-free, efficient, and interactive platform for everyday tasks. Its core strength lies in the seamless combination of visual and auditory inputs, creating a highly immersive and accessible technological solution that empowers users to engage with their devices in innovative ways.

The potential applications of IntelleX span across a variety of industries, each benefiting from its unique combination of AI-driven features and multi-modal interaction. For instance:

1. **Education**
   IntelleX can transform the educational sector by providing students with interactive learning experiences. Through gesture and voice-based controls, students can explore complex concepts, such as scientific phenomena, historical events, and geographical features, in a more dynamic and engaging manner.
2. **Healthcare and Accessibility**
   In healthcare, IntelleX can be a powerful assistive tool for people with disabilities, enabling them to perform tasks independently through voice commands and gestures. The object detection feature can help those with visual impairments by providing real-time feedback and aiding in navigation and object identification.
3. **Smart Environments**
   IntelleX is ideal for use in smart homes and cities, where real-time object tracking and intelligent interaction can optimize energy consumption, enhance security through surveillance, and facilitate smooth communication between various smart devices and systems.
4. **Entertainment and Gaming**
   In the entertainment sector, IntelleX opens new possibilities for immersive AR experiences, where real-world objects can be tracked and recognized in real time. This technology can create innovative gaming and interactive storytelling experiences, merging physical and digital worlds.
5. **Retail and E-commerce**
   IntelleX can enhance the retail experience by enabling virtual try-ons or offering dynamic product information through augmented reality. This capability can drive customer engagement, help make more informed purchasing decisions, and enhance overall satisfaction.

As the project develops, it continues to bridge the gap between human interaction and machine intelligence, enhancing the user experience. However, there remain areas for improvement, particularly in refining gesture recognition accuracy, enhancing the speed and adaptability of object detection, and ensuring system scalability and performance for diverse applications.

In conclusion, IntelleX is not just an innovative tool; it is poised to transform the way we interact with both the digital and physical worlds. With its combination of cutting-edge AI, real-time object detection, and multi-modal interaction, IntelleX has the potential to revolutionize industries such as education, healthcare, smart environments, entertainment, and retail. As the technology matures, it will continue to offer new possibilities, making technology more accessible, intuitive, and impactful for users across the globe.

# 7. Future Scope

As the demand for intelligent and accessible digital assistants continues to grow, **IntelleX** offers substantial opportunities for advancement across various domains. Its multi-modal framework, combining gesture, voice, and vision, can be scaled and enhanced to support broader functionalities and user needs.

1. **Advanced Gesture Recognition**
   Future updates can integrate deep learning models for more accurate gesture tracking, recognizing subtle movements, multi-finger control, and dynamic hand gestures. This would increase responsiveness and adaptability in diverse environments.
2. **IoT and Smart Device Integration**
   IntelleX can be integrated with smart home and IoT devices to enable seamless control through gestures and voice. Users could manage lights, thermostats, appliances, and security systems, making homes and offices more intelligent and automated.
3. **Augmented Reality and Object Interaction**
   Expanding into AR, IntelleX could allow users to interact with real-world objects using visual overlays. This could enhance experiences in education, navigation, and gaming by offering real-time guidance and contextual information.
4. **Multi-Language & Accent Support**
   To ensure global usability, future versions of IntelleX could support multilingual voice commands and adapt to different accents using advanced natural language processing (NLP) models.
5. **Accessibility Expansion**
   IntelleX can play a transformative role for users with disabilities by integrating sign language detection, voice-guided assistance, and real-time object descriptions, offering improved independence and usability.
6. **Cloud and Edge Computing**
   Utilizing edge and cloud computing can reduce latency and handle more complex computations, enabling faster processing, real-time updates, and scalability across multiple platforms and devices.
7. **Personalization and Predictive AI**
   With machine learning, IntelleX can learn user behavior, provide context-aware recommendations, and anticipate actions based on previous interactions, making the assistant more intelligent and helpful over time.
8. **AR/VR Integration for Virtual Workspaces**
   Future developments may bring IntelleX into AR/VR environments, enabling immersive interaction for online education, virtual meetings, and collaborative workspaces through gestures and voice.
9. **Security and Privacy Enhancements**
   As AI assistants handle sensitive data, implementing secure authentication, encrypted communication, and user consent mechanisms will be essential for maintaining trust and protecting user privacy.

**In conclusion**, the future of IntelleX is filled with opportunities for technological and societal impact. By continuously integrating emerging technologies, it can evolve into a smarter, more inclusive, and highly adaptable assistant that enhances productivity, accessibility, and human-computer interaction in meaningful ways.

# 8.References

- Python **Software Foundation** – https://www.python.org

- Flask **Documentation** – https://flask.palletsprojects.com

- OpenCV **Library** – https://opencv.org

- Mediapipe **by Google** – https://mediapipe.dev

- Speech **Recognition Library** – https://pypi.org/project/SpeechRecognition/

- pyttsx**3 Text-to-Speech** – https://pypi.org/project/pyttsx3/

- Pycaw **– Python Core Audio Windows Library** – https://github.com/AndreMiras/pycaw

- TensorFlow **Documentation** – https://www.tensorflow.org

- Keras **API Reference** – https://keras.io

- NumPy **Library** – https://numpy.org

- Matplotlib **Documentation** – https://matplotlib.org

- Jupyter **Notebook Project** – https://jupyter.org

- Zhang, X., Zhao, J., & LeCun, Y. (2015). *Character-level Convolutional Networks for Text Classification*. arXiv preprint arXiv:1509.01626.

- He, K., Gkioxari, G., Dollár, P., & Girshick, R. (2017). *Mask R-CNN*. Proceedings of the IEEE International Conference on Computer Vision (ICCV).

# 9. Appendix

- **Code Snippets:**

  - Gesture Recognition Module: *Utilizes Mediapipe and OpenCV to detect and interpret hand gestures in real-time for tasks like media control, navigation, and system operations.*
  - Voice Command Module: *Employs the SpeechRecognition library to capture and process voice inputs, paired with pyttsx3 for delivering spoken responses, enhancing interactivity.*
- **User Interface Design:**
  - The user interface includes a live webcam feed for gesture tracking and a dynamic voice input console that displays recognized commands.
  - Additional features include settings customization, a help section, and toggle options for enabling/disabling specific modules such as object detection and audio feedback.
- **System Operation Flowchart:**
  - The system starts with capturing input (gesture or voice) from the user.
  - The backend (Flask) processes this input, matches it with a corresponding function, and executes the required system or application task.
  - The output is then delivered through visual feedback or audio responses.
- **System Requirements:**
  - **Hardware:** Functional webcam and microphone, 8GB RAM or higher, Intel Core i5 processor or equivalent for smooth real-time processing.
  - **Software:** Python 3.8+, required libraries (OpenCV, Mediapipe, SpeechRecognition, pyttsx3, Pycaw, TensorFlow/Keras), and OS support for Windows/Linux/Mac.
- **Testing and Evaluation:**
  - The gesture recognition module was tested across varied lighting conditions and hand shapes to ensure robustness.
  - Voice recognition accuracy was validated in both quiet and noisy environments to evaluate its reliability in real-life scenarios.
- **Future Enhancements:**
  - Integration of advanced real-time object detection using pre-trained DNNs.
  - Porting the system to mobile and wearable platforms for greater accessibility.
  - Addition of language support for multilingual voice commands and personalized user profiles.