

PransFly

Ponnekanti Pranathi^{1,2*}

^{1*} Computer Science and Engineering, IIT Goa, Farmagudi, Ponda,
403401, Goa, India.

Corresponding author(s). E-mail(s):
ponnekanti.pranathi.21031@iitgoa.ac.in;

Abstract

This project explores the integration of advanced computer vision, natural language processing, and drone control technologies to enhance autonomous aerial systems. Leveraging the YOLO object detection model for real-time visual processing and OpenAI’s GPT-based Large Language Models (LLMs) for contextual reasoning, this system enables drones to perform complex tasks guided by human-like comprehension of instructions.

Keywords: Computer Vision, Large Language Models, Drone, YOLO Object Detection

1 Introduction

Controlling drones through natural language brings a new level of intuitiveness to human-drone interaction, simplifying operations while opening doors for advanced AI systems to autonomously manage drone activities. This project leverages the capabilities of large language models (LLMs) to translate task instructions in natural language into executable commands for drones, revolutionizing how drones can be operated. By integrating a remote LLM (GPT-4o-mini) with cutting-edge object detection (YOLOv10x), facial recognition (face recognition library), and speech recognition (e.g., WebKit Speech Recognition), this project demonstrates an innovative and cohesive approach to drone control.

The system employs a multi-step process where user queries and real-time image frames from the drone’s camera are sent to the remote LLM. The LLM processes these inputs—combining object detection results and user prompts—to generate actionable responses. Based on the response, the system commands the drone to perform specific tasks, such as detecting, identifying, or navigating around objects or individuals. The

incorporation of face recognition ensures user-specific tasks, such as identifying authorized personnel, while speech recognition provides an additional layer of accessibility, allowing for voice-based commands.

This integration results in a complex yet highly effective system capable of intelligent decision-making in real-world scenarios. The project's architecture not only bridges natural language understanding and robotics but also sets a foundation for the future of language-driven automation in drone technology.

2 Results

2.1 Object Detection and Classification

The system successfully uses the YOLOv10x model to detect and classify objects in the video feed in real-time. Detected objects are stored in a list (detected objects) and annotated on the video feed with bounding boxes and labels. The bounding boxes and annotations are dynamically updated on the video stream to reflect the current detections.

2.2 "Detect Only" Mode

Users can specify a target object for detection by including the phrase "detect only [target]" in their query. When "detect only" mode is active, the application focuses solely on detecting the specified target object. If the target object is detected, it is highlighted with a bounding box and label on the video stream.

2.3 Facial Recognition

For objects classified as "person," the system attempts to identify faces using the face recognition library. Recognized faces are matched against a database of known encodings (known face encodings), and the associated names (known face names) are displayed. For unknown faces, the user is prompted to provide a name for the newly detected face, which can then be added to the database.

2.4 Rotational Scanning

If an object relevant to the user's query is not immediately found, the system performs a virtual rotation (incrementing by 45 degrees up to 360 degrees) to scan for the target object. After a full rotation without finding the target object, the system concludes the task and notifies the user.

2.5 Integration with OpenAI GPT

The system generates a prompt for OpenAI's GPT model, combining the detected objects and the user's query, and includes a base64-encoded image of the current frame for contextual understanding. GPT provides a response indicating whether it can answer the query based on the detected objects, categorized as "yes" or "no." The application's behavior (e.g., stopping the task or continuing detection) is guided by GPT's response.

2.6 Real-Time Feedback via WebSockets

The results of detections, including annotated video frames and responses to user queries, are delivered in real time through WebSocket events. The client interface receives continuous updates of the detection progress and task status.

2.7 Multimodal Interaction

The system supports both textual queries and image data, enabling seamless interaction between video feed analysis and natural language understanding.

2.8 Voice command integration

Used Webspeech API, a speech recognition API which supports real time voice-to-text transcriptions.

3 Workflow of the application

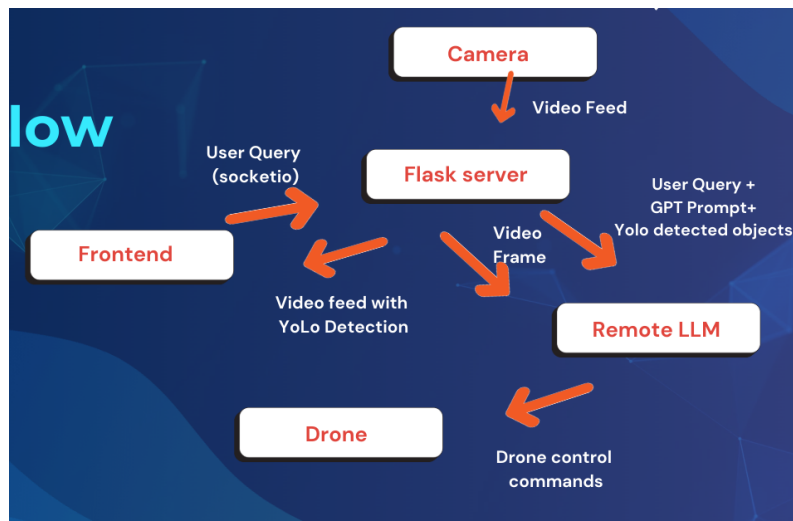


Fig. 1 Workflow of the application.

4 Figures

These are some of the screenshots of my application and its features. It has a user friendly interface to display the drone PoV(Drone Point of View) and a chat interface which is connected to the OpenAI's GPT model. It has a mic to enable speech recognition which does the same functionality as the text input.

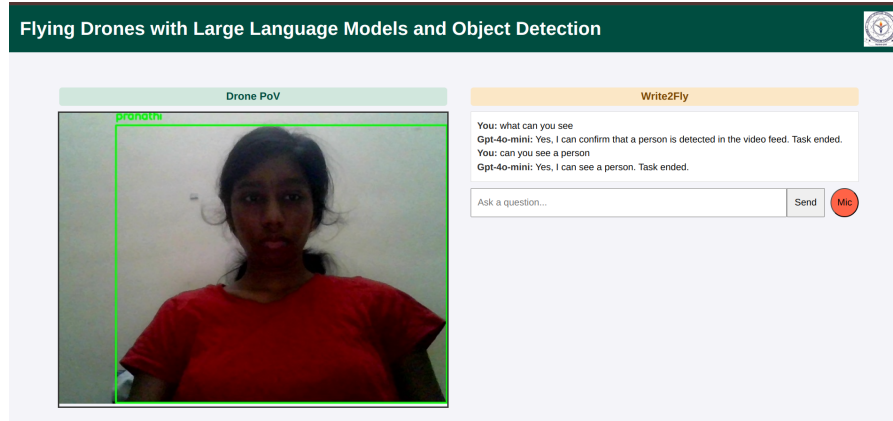


Fig. 2 UserInterface of PransFly

5 Algorithms, Program codes and Listings

5.1 Algorithms and Models used

Utilized yolov10x for object detection and manipulation of yolo focus by implementing detect only functionality. I have used 'face recognition' library and functionalities to identify and store the faces. Utilized socketio library for connection between the frontend and backend flask server(to send the queries and receive the response).

5.2 Program codes and Listings

5.2.1 Yolo Detection and Frame Annotation

The following listing demonstrates how YOLO is integrated for object detection and frame annotation:

Algorithm 1 YOLO Detection and Annotation

```

1: Results  $\leftarrow$  model.predict(source=frame, show=False, verbose=False)
2: Detected_objects  $\leftarrow$  [model.names[int(cls)] for cls in results[0].boxes.cls]
3: for all box, cls, conf in zip(results[0].boxes.xyxy, results[0].boxes.cls,
   results[0].boxes.conf) do
4:   obj_name  $\leftarrow$  model.names[int(cls)]
5:    $x_1, y_1, x_2, y_2 \leftarrow$  map(int, box)
6:   Draw rectangle with  $(x_1, y_1), (x_2, y_2)$  and color (0, 0, 255)
7:   Add label obj_name at  $(x_1, y_1 - 10)$ 
8: end for

```

Annotation Code:

```

1 # Predict results using YOLO model
2 results = model.predict(source=frame, show=False, verbose=False)
3
4 # Extract detected object names
5 detected_objects = [model.names[int(cls)] for cls in results[0].boxes.cls]
6
7 # Loop through detected bounding boxes, classes, and confidences
8 for box, cls, conf in zip(results[0].boxes.xyxy, results[0].boxes.cls, results[0].
   boxes.conf):
9     obj_name = model.names[int(cls)]
10    x1, y1, x2, y2 = map(int, box)
11
12    # Draw rectangle around the detected object
13    cv2.rectangle(frame, (x1, y1), (x2, y2), (0, 0, 255), 2)
14
15    # Add label above the bounding box
16    cv2.putText(frame, obj_name, (x1, y1 - 10), cv2.FONT_HERSHEY_SIMPLEX, 0.6, (0,
   0, 255), 2)

```

5.2.2 Face Recognition Integration

This listing explains how face recognition is used to match or register new faces:

Algorithm 2 Face Recognition Logic for YOLO Detected Person

```

1: if obj_name == "person" then
2:   cropped_face ← frame[y1 : y2, x1 : x2]
3:   rgb_face ← Convert(cropped_face, BGR → RGB)
4:   face_encodings ← FaceEncodings(rgb_face)
5:   if face_encodings ≠ ∅ then
6:     matches ← CompareFaces(known_face_encodings, face_encodings[0], tolerance =
       0.6)
7:     if True ∈ matches then
8:       match_index ← IndexOf(True, matches)
9:       iname ← known_face_names[match_index]
10:    else
11:      name ← Input("Name : ")
12:      if name ≠ ∅ and name ∉ known_face_names then
13:        Append(known_face_encodings, face_encodings[0])
14:        Append(known_face_names, name)
15:      end if
16:    end if
17:  end if
18: end if

```

```

1 if obj_name == "person":
2     # Crop and preprocess the detected face
3     cropped_face = frame[y1:y2, x1:x2]
4     rgb_face = cv2.cvtColor(cropped_face, cv2.COLOR_BGR2RGB)
5
6     # Generate face encodings
7     face_encodings = face_recognition.face_encodings(rgb_face)

```

```

8
9     if face_encodings:
10         # Compare face encoding with known encodings
11         matches = face_recognition.compare_faces(known_face_encodings,
12             face_encodings[0], tolerance=0.6)
13
14         if True in matches:
15             # Match found, retrieve the name
16             match_index = matches.index(True)
17             iname = known_face_names[match_index]
18         else:
19             # New face, prompt for a name
20             print("New face detected. Enter the person's name:")
21             name = input("Name: ").strip()
22
23             if name and name not in known_face_names:
24                 # Save new face encoding and name
25                 known_face_encodings.append(face_encodings[0])
26                 known_face_names.append(name)

```

5.2.3 Flask Web Application and WebSocket Integration

This listing shows the code structure for creating the video feed and handling client-server interactions:

Algorithm 3 Video Streaming and Question Handling via Flask and SocketIO

```

1: function VIDEO_STREAM
2:     return Response(generate_frames(), mimetype =
        'multipart/x-mixed-replace; boundary=frame')
3: end function
4: function HANDLE_QUESTION(data)
5:     global detected_objects, last_frame
6:     question ← data['question']
7:     response_message ← query_openai(gpt_prompt, encode_image(last_frame))
8:     emit('receive_answer', {'response': response_message})
9: end function

```

Flask and socket io integration Code:

```

1 from flask import Flask, Response, render_template
2 from flask_socketio import SocketIO
3 import cv2
4 import base64
5
6 app = Flask(__name__)
7 socketio = SocketIO(app)
8
9 # Video streaming route
10 @app.route('/video_feed')
11 def video_stream():
12     return Response(generate_frames(), mimetype='multipart/x-mixed-replace;
13         boundary=frame')
14
15 # Route to render HTML page
16 @app.route('/video_page')
17 def video_page():
18     return render_template('index.html') # Render the HTML template

```

```

18
19 # Function to encode the image frame
20 def encode_image(image):
21     _, buffer = cv2.imencode('.jpg', image)
22     return base64.b64encode(buffer).decode('utf-8')
23
24 # WebSocket event for processing questions from the client
25 @socketio.on('send_question')
26 def handle_question(data):
27     global detected_objects, task_ended, current_angle, last_frame
28
29     question = data['question']
30     response_message = ""
31
32     # Example logic for processing the question
33     if question == "Is the task done?":
34         response_message = "Yes, the task is complete."
35
36     # Emit the response back to the client
37     emit('receive_answer', {'response': f"{response_message} Task ended."})
38
39 if __name__ == '__main__':
40     socketio.run(app)

```

5.2.4 Integration with OpenAI GPT

I have used the model 'gpt-4o-mini' through openai api. This function takes an image along with the user's query and gpt prompt as input and sends the query to gpt model through the api. The response that it receives back is again emitted using socketio.

Algorithm 4 Query OpenAI GPT-4 with Base64 Image Input

```

1: function QUERY_OPENAI(gpt_prompt, base64_image)
2:   Set openai.api_key ← YOUR_API_KEY
3:   Create response with OpenAI API:
4:     Model: "gpt-4o-mini"
5:     Messages:
6:       {"role": "user", "content": gpt_prompt}
7:       {"role": "user", "image_url": "data:image/jpeg;base64," + base64_image}
8:   return response.choices[0].message.content
9: end function

```

Code for Querying OpenAI GPT-4 with Base64 Image Input:

```

1 import openai
2
3 def query_openai(gpt_prompt, base64_image):
4     openai.api_key = "YOUR_API_KEY" # Replace with your actual API key
5     response = openai.chat.completions.create(
6         model="gpt-4o-mini",
7         messages=[
8             {
9                 "role": "user",
10                "content": gpt_prompt,
11            },
12            {
13                "role": "user",

```

```

14         "image_url": {
15             "url": f"data:image/jpeg;base64,{base64_image}"
16         },
17     },
18 ]
19 )
20 return response.choices[0].message.content

```

5.2.5 Task completion and rotation

This code implements a task loop that uses object detection and OpenAI's API to determine the relevance of detected objects in a video feed to a user's question. The loop iteratively rotates a camera or view by 45 degrees, checking for relevant objects within each frame until either a relevant object is found or a full 360-degree rotation is completed. For each frame, if objects are detected and a previous frame is available, the last frame is encoded in Base64 format, and a prompt is generated for OpenAI. The prompt provides details about the detected objects and asks OpenAI to determine whether the objects are relevant to the user's query, responding with either "yes" or "no." Based on OpenAI's response, the task may end if relevance is detected ("yes"), or the camera rotation continues if not ("no"). If the camera completes a full rotation without finding relevant objects, the task ends with a message indicating that no relevant objects were detected. Throughout the process, feedback is printed to the console, and the final response is sent to a client using 'emit'. The below screenshot shows how the commands are being printed on the console to rotate or navigate whenever there is no positive relevance found.

```

No positive relevance in response. Command: Rotate to 45 degrees
No, I cannot determine if there is a cell phone present in the image.
No positive relevance in response. Command: Rotate to 90 degrees
No, I cannot determine if there is a cell phone visible in the image.
No positive relevance in response. Command: Rotate to 135 degrees
No, I cannot determine if there is a cell phone visible in the image.
No positive relevance in response. Command: Rotate to 180 degrees
No, I cannot determine if a cell phone is visible based on the current information.
No positive relevance in response. Command: Rotate to 225 degrees
No, I cannot confirm if there is a cell phone present in the image.
No positive relevance in response. Command: Rotate to 270 degrees
No, I cannot confirm whether a cell phone is present based on the information provided.
No positive relevance in response. Command: Rotate to 315 degrees
No, I cannot confirm the presence of a cell phone based on the current information.
No positive relevance in response. Command: Rotate to 360 degrees
Completed a full rotation without detecting the relevant object.

```

Fig. 3 Commands to navigate in console

Code for Task Execution with OpenAI and Object Detection:

```

1 import openai
2 import socketio
3 import cv2
4 import base64
5

```

Algorithm 5 Task Execution with OpenAI and Object Detection

```
1: Initialize task_ended as False, current_angle as 0
2: while not task_ended and current_angle < 360 do
3:   if detected_objects and last_frame  $\neq$  None then
4:     base64_image  $\leftarrow$  Encode(last_frame)
5:     gpt_prompt  $\leftarrow$  "The user asked: 'question'.
      The following objects are currently detected in the video feed: detected_objects.
      Based on this information, add 'yes' with a response if relevant, otherwise 'no'."
6:     response_message  $\leftarrow$  QueryOpenAI(gpt_prompt, base64_image)
7:     Print(response_message)
8:     if "yes" in response_message.lower() then
9:       Print("Positive response detected. Ending task.")
10:      task_ended  $\leftarrow$  True
11:      break
12:     else if "no" in response_message.lower() then
13:       Print("No positive relevance. Rotate to current_angle + 45 degrees.")
14:     end if
15:   end if
16:   current_angle  $\leftarrow$  current_angle + 45
17:   socketio.sleep(10)
18: end while
19: if current_angle  $\geq$  360 then
20:   task_ended  $\leftarrow$  True
21:   Print("Completed a full rotation without detecting the relevant object.")
22:   response_message  $\leftarrow$  "No relevant object found after full rotation."
23: end if
24: Emit("receive_answer", {'response': response_message + " Task ended."})
```

```
6 # Initialize global variables
7 task_ended = False
8 current_angle = 0
9
10 def encode_image(image):
11     _, buffer = cv2.imencode('.jpg', image)
12     return base64.b64encode(buffer).decode('utf-8')
13
14 def query_openai(gpt_prompt, base64_image):
15     openai.api_key = "YOUR_API_KEY" # Replace with your actual API key
16     response = openai.chat.completions.create(
17         model="gpt-4o-mini",
18         messages=[
19             {
20                 "role": "user",
21                 "content": gpt_prompt,
22             },
23             {
24                 "role": "user",
25                 "image_url": {
26                     "url": f"data:image/jpeg;base64,{base64_image}"
27                 },
28             },
29         ]
30     )
```

```

31     return response.choices[0].message.content
32
33 while not task_ended and current_angle < 360:
34     if detected_objects and last_frame is not None:
35         base64_image = encode_image(last_frame)
36         gpt_prompt = f"The user asked: '{question}'. \n\nThe following objects are
            currently detected in the video feed: {detected_objects}. Based on
            this information, add 'yes' with a response if relevant, otherwise 'no'
            ,."
37         response_message = query_openai(gpt_prompt, base64_image)
38         print(response_message)
39
40         if "yes" in response_message.lower():
41             print("Positive response detected. Ending task.")
42             task_ended = True
43             break
44         elif "no" in response_message.lower():
45             print(f"No positive relevance. Rotate to {current_angle + 45} degrees."
                )
46
47         current_angle += 45
48         socketio.sleep(10)
49
50 if current_angle >= 360:
51     task_ended = True
52     print("Completed a full rotation without detecting the relevant object.")
53     response_message = "No relevant object found after full rotation."
54
55 socketio.emit("receive_answer", {'response': response_message + " Task ended."})

```

5.2.6 Voice command integration

Algorithm 6 Speech Recognition Integration

```

1: function STARTLISTENING
2:     Initialize speech recognition system
3:     Set language to "en-US"
4:     Set recognition mode to "interimResults"
5:     Start recognition
6:     if Recognition started then
7:         Print "Voice recognition started"
8:     end if
9:     while Recognition continues do
10:         Capture interim results
11:         Update input field with transcript
12:     end while
13:     if Recognition ended then
14:         Print "Voice recognition ended"
15:         Send the transcribed text as a message
16:     end if
17:     if Error occurs then
18:         Print "Speech recognition error"
19:     end if
20: end function

```

Below is the JavaScript code used for implementing speech recognition:

Listing 1 Speech Recognition Integration Code

```
const recognition = new (window.SpeechRecognition)();
recognition.lang = 'en-US';
recognition.interimResults = true;

// Start listening on microphone button click
function startListening() {
    recognition.start();
    recognition.onstart = function() {
        console.log('Voice recognition started');
    };

    recognition.onresult = function(event) {
        let transcript = '';
        for (let i = event.resultIndex; i < event.results.length; i++) {
            transcript += event.results[i][0].transcript;
        }

        // Set the transcript in the input field
        document.getElementById("userInput").value = transcript;
    };

    recognition.onerror = function(event) {
        console.error('Speech recognition error', event);
    };

    recognition.onend = function() {
        console.log('Voice recognition ended');
        sendMessage(); // Send the transcribed text after recognition ends
    };
}
```

5.2.7 GPT-4's sentiment analysis

Making use of powerful gpt model's sentiment analysis functionality to check whether the relevance is found or not in the response. This step is useful in many other cases when we have to check relevance of the response to the visual representations and perform the navigation of the drone accordingly.

Sentiment Analysis

```
1 gpt_prompt = (f"The user asked: '{question}'. The following objects are currently
2     detected in the video feed: "
3         f"{'', '.join(detected_objects)}. Based on this information, please
4             add a 'yes' along with the gpt response if it is possible to
5             provide an answer relevant to the user's question, "
6             "or 'no' along with the gpt response if it is not possible.")
7
8 response_message = query_openai(gpt_prompt, base64_image)
9 print(response_message) # Print the response to the console
10
11 if "yes" in response_message.lower(): # Positive sentiment suggests relevance
12     print("Positive response detected. Ending task.")
13     task_ended = True
14     break
15 elif "no" in response_message.lower():
16     print(f"No positive relevance in response. Command: Rotate to {current_angle +
17         45} degrees")
```

Algorithm 7 Sentiment Analysis

```
1: function QUERY_GPT_RESPONSE(question, detected_objects, base64_image)
2:   Generate gpt_prompt with user question and detected objects
3:   Set the prompt to ask if it is possible to provide an answer relevant to the
   question
4:   Send the gpt_prompt along with base64_image to OpenAI API
5:   Capture the response in response_message
6:   Print the response_message
7:   if "yes" in response_message then
8:     Print "Positive response detected. Ending task."
9:     Set task_ended  $\leftarrow$  True
10:    Exit the loop
11:  else if "no" in response_message then
12:    Print "No positive relevance in response."
13:    Print "Command: Rotate to current_angle + 45 degrees"
14:  end if
15: end function
```

6 Discussion

Integrating drones with large language models (LLMs) like GPT unlocks a wide range of advanced applications, enabling real-time analysis and contextual decision-making. For security, drones equipped with object detection can identify weapons, suspicious activities, or intruders, providing natural language reports about threats and understanding risks in scenarios like crowded areas or emergencies. In traffic monitoring, drones can recognize car models, brands, and modifications, cross-reference data with stolen vehicle databases, and answer natural language queries like, "How many red cars passed this street today?" In agriculture, drones can analyze crop types, recommend fertilizers or pest control solutions, and respond to farmer queries such as, "Which area needs more irrigation?" They can also explain environmental conditions impacting yield, like pest infestations. For healthcare, drones can recognize emergencies, explain medical kit contents, and offer step-by-step CPR instructions or locate the nearest hospital. These innovations demonstrate the transformative potential of combining drones with LLMs, enabling applications far beyond traditional systems and redefining efficiency and intelligence across multiple industries.

7 Conclusion

This project demonstrates a simulation of drone-like behavior, utilizing advanced functionalities implemented through YOLO-based object detection and the integration of large language models (LLMs). The code supports a "detect only" mode, enabling the system to prioritize and follow a specific target object in real-time, effectively mimicking a drone's ability to focus on a subject. Additionally, it incorporates a task rotation mechanism, whereby the system incrementally rotates by 45 degrees when no positive relevance is detected, systematically scanning the environment for the desired object.

The simulation concludes tasks upon either detecting the target object or completing a full 360-degree rotation. Given the unavailability of a physical drone, this implementation strives to simulate the core aspects of such a system, showcasing its potential to enhance drone operations through intelligent object detection and decision-making powered by LLMs.

References

- [1] Redmon, Joseph, Divvala, Santosh, Girshick, Ross, and Farhadi, Ali. *You Only Look Once: Unified, Real-Time Object Detection*. arXiv preprint arXiv:1506.02640, 2015.
- [2] Ultralytics. *Ultralytics YOLO*. Available at: <https://github.com/ultralytics/ultralytics>.
- [3] OpenAI. *OpenAI API*. Available at: <https://platform.openai.com/docs/>.
- [4] Ageitgey, Adam. *Face Recognition*. Available at: https://github.com/ageitgey/face_recognition.
- [5] Pallets Projects. *Flask: The Python microframework for building web applications*. Available at: <https://flask.palletsprojects.com/>.
- [6] Grinberg, Miguel. *Flask-SocketIO*. Available at: <https://flask-socketio.readthedocs.io/>.
- [7] Bradski, Gary and Kaehler, Adrian. *OpenCV: Open Source Computer Vision Library*. Available at: <https://opencv.org/>.
- [8] Python Software Foundation. *Base64 Encoding and Decoding in Python*. Available at: <https://docs.python.org/3/library/base64.html>.