

Object Tracking Algorithms

ALHassan Baligh Hassan ALShami

Abdulghani Khaled Abdul Al-Wajeih

Mohammed Khaled Al-Alaimi

Osama Saif Abdo Rafea

Supervisors

Dr. Amin Mohamed Ahsan Shayea

Dr. Hisham Aqlan

**FACULTY OF COMPUTER SCIENCE AND INFORMATION
TECHNOLOGY INTERNATIONAL UNIVERSITY OF TECHNOLOGY
TWINTech 2024 – 2025**

Object Tracking Algorithms

ALHassan Baligh Hassan ALShami
Abdulghani Khaled Abdul Al-Wajeeh
Mohammed Khaled Al-Alaimi
Osama Saif Abdo Rafea

A report submitted in fulfillment of the requirements for the award of the degree of
Bachelor of Science (Computer Science)

FACULTY OF COMPUTER SCIENCE AND INFORMATION
TECHNOLOGY INTERNATIONAL UNIVERSITY OF
TECHNOLOGY TWINTech

June, 2025

INTERNATION UNIVERSITY OF TECHNOLOGY TWINTech

DECLARATION OF PROJECT COPYRIGHT

Author's full name: ALHassan Baligh Hassan ALShami

Date of Birth : 2002 / 10 / 01

Author's full name: Abdulghani Khaled Al-Wajeeh

Date of Birth : 1998 / 00 / 00

Author's full name: Mohammed Khaled Al-Alaimi

Date of Birth : 2000 / 00 / 00

Author's full name: Osama Saif Abdo Rafea

Date of Birth : 2003 / 00 / 00

Title : Object Tracking Algorithms

Academic Session : 2024/ 2025 (....)

I declare that this Dissertation is classified as:

☐

CONFIDENTIAL

(Official Letter)

☐

RESTRICTED

(Official Letter)

☐

OPEN ACCESS

I/We agree that my Report to be published as
online open access (full text)

I/We acknowledged that International University of Technology Twin tech reserves the right as follows:

1. The Report is the property of International University of Technology Twintech
2. The Library of International University of Technology Twintech has the right to make copies for the purpose of research or academic exchange.

Certified by:

.....
SIGNATURE OF STUDENTS

06143880

IC No/PASSPORT

Date:

.....
SIGNATURE OF SUPERVISORS

Dr. Amin Mohamed Ahsan Shaye
, Dr. Hisham Aqlan

NAME OF SUPERVISORS

Date:

I hereby declare that I have read this Report titled “Object Tracking Algorithms” and, in my opinion, is sufficient in term of scope and quality for the award of the degree of Bachelor Science of Computer Science.

Signature:

Name of Supervisor: Amin Mohamed Ahsan Shayea

Date: June, 2025

Signature:

Name of Supervisor: Hisham Aqlan

Date: June, 2025

We hereby declare that this PROJECT entitled “Object Tracking Algorithms” is the result of our own research except as cited in the references. The report has not been accepted for any degree and is not concurrently submitted in candidature of any other degree.

Signature:

Name: ALHassan Baligh Hassan ALShami

Signature:

Name: Abdulghani Khaled Al-Wajeeh

Signature:

Name: Mohammed Khaled Al-Alaimi

Signature:

Name: Osama Saif Abdo Rafea

Date:

DEDICATION

This project is dedicated to our precious parents who have been a constant source of support and encouragement during life. We also dedicated this work to Dr. Amin Shayea who guided and helped us through the whole process of this project and kept us on track until we successfully complete it.

ACKNOWLEDGMENTS

“In the Name of Allah and peace and blessings of Allah be upon our prophet, Mohammad, and upon all his family and companions”. In the first place, words cannot express the gratitude to Allah who gave us the inspiration and determination to complete Object Tracking Algorithms Project successfully. The team could not have undertaken this journey without the families of each member who gave us the support and encouragement to achieve our desired goal which is building this project at the highest level we could. Our project team extremely grateful to the supportive supervisor, Dr. Amin Shayea for his significant instructions and guidance. We would be remiss in not mentioning our dean Dr. Hamzah Jamel for the key information that we have taken through several courses related to software analysis and system design.

Abstract

This project presents Object Tracking Algorithms an AI-powered visual perception system designed to simulate real-time object detection, multi-target tracking, depth estimation, and gesture-based interaction using advanced computer vision techniques. At its core, the system employs a high-resolution camera and integrates state-of-the-art models including YOLOv8 for object detection, ByteTrack for identity-preserving multi-object tracking, MiDaS for monocular depth estimation, and MediaPipe for gesture and human pose recognition. Unlike traditional hardware-dependent systems, this project simulates the "eyes" and logic of an intelligent agent capable of autonomously following objects, estimating distances, and responding to visual commands. The tracking system combines image-based depth inference with dynamic control logic based on a modular PID controller, further refined by Kalman filtering for smoother transitions. Gesture-based commands, such as "FIST" to reload parameters or "POINT" to reset tracking, allow intuitive interaction with the system. Depth fusion with optional LIDAR enhances obstacle detection, and visual alerts are triggered when obstructions interrupt the line of sight to the tracked object. Designed for simulation and real-time visualization, the architecture is modular, GPU-accelerated, and adaptable for deployment on embedded platforms like Raspberry Pi or NVIDIA Jetson. The system offers a robust foundation for intelligent surveillance, educational robotics, smart camera systems, and experimental human-AI interaction.

.

List of Content

DECLARATION OF PROJECT COPYRIGHT	3
DEDICATION	6
ACKNOWLEDGMENTS	7
Abstract	8
Chapter 1	13
Introduction	13
1.1 An Overview:	13
1.2 Problem Background:	14
1.3 Problem Statement:	14
1.4 Research Questions:	15
1.5 Research Aim:	15
1.6 Research Objectives:	15
1.7 Scope:	16
1.8 Outline for Project	16
1.9 Terms Definition	16
1.10 Summary of Research Contributions	17
1.11 Project Outline	17
Chapter 2	18
Literature Review	18
2.1 Introduction:	18
2.2 Visual Sensing	18
2.2.1 Image Processing and Computer Vision Algorithms	19
2.2.3 Navigation Integration and Control Logic	19
2.2.4 System Challenges	20
2.3 Obstacle Awareness and Avoidance Logic	20
2.3.1 Object Recognition	20
2.3.2 Depth Perception	20
2.3.3 Motion Smoothing with Kalman	20
2.3.4 Avoidance Strategy	21
2.3.5 Role of AI and ML Models	21
2.3.6 Real-Time Performance	21
2.3.7 Common Challenges	21
2.3.8 Applications	22
2.4 Datasets and Feature Sources	23

2.4.1 Reference Datasets	23
2.4.2 Pretraining Datasets	23
2.4.3 Synthetic Simulation Datasets.....	23
2.4.4 Dataset Challenges	23
2.5 Performance Evaluation.....	24
2.5.1 Streaming Pipeline	24
2.6 Related Work.....	24
2.6.1 Navigation Models	24
2.6.2 Obstacle Avoidance Models	24
2.6.3 Research Areas	24
2.6.4 Future Opportunities	25
2.7 Summary	26
Chapter 3	27
Research Methodology	27
3.1 Introduction	27
3.2 Study Framework	28
3.2.1 Visual Perception.....	29
3.2.2 Gesture Recognition & Pose Estimation	30
3.2.3 Control Simulation	31
3.2.4 Logging & Visualization	31
3.3 Dataset.....	31
3.3.1 Public Datasets (Indirect Use).....	31
3.3.2 Custom Simulation Scenarios	32
3.4 Proposed Methods.....	33
3.4.1 Object Detection Framework.....	33
3.4.2 Depth Estimation and Obstacle Awareness	34
3.4.3 Real-Time Processing Pipeline.....	35
3.5 Experimental Design	36
3.6 Performance Evaluation.....	37
3.6.1 Performance Metrics	38
3.6.2 Sample Results	38
3.7 Summary	39
Chapter 4	40
System Design and Implementation	40
4.1 Introduction:	40
4.2 Methods	41
4.2.1 Motivation	41
4.2.2 System Architecture	42

4.2.3 Implementation Steps	46
4.2.4 Real-Time Video Processing	46
4.3 Introduction:	47
4.4 Summary	48
Chapter 5	49
References	53

Chapter 1

Introduction

1.1 An Overview:

Object tracking plays a crucial role in advancing the autonomy of Unmanned Aerial Vehicles (UAVs), commonly referred to as drones. These systems are transforming fields such as surveillance, inspection, and robotics by enabling drones to visually follow and interpret dynamic targets. With the evolution of AI and computer vision, drones are now capable of tracking multiple objects, estimating their depth, and responding to human interaction in real time—without the need for GPS or predefined paths. This project, titled Object Tracking Algorithms, presents a modular simulation of such a system. It focuses on object detection, multi-target tracking, gesture interpretation, and obstacle awareness using only a camera input. By integrating advanced models such as YOLOv8, ByteTrack, MiDaS, and MediaPipe, the system demonstrates intelligent, visual-based control logic with depth-aware tracking and gesture-driven interaction.

1.2 Problem Background:

Traditional drone navigation methods that rely on GPS or waypoint planning often fail in indoor or cluttered environments. Visual tracking systems provide a flexible alternative by analyzing image streams to detect and follow targets. However, building such systems introduces several challenges:

- High computational demand, especially for real-time tracking and gesture analysis
- Errors caused by occlusions, dynamic lighting, or fast-moving targets
- Limited ability to estimate object depth and avoid obstacles
- Integration complexity when fusing detection, tracking, and pose estimation

Object Tracking Algorithms addresses these issues by combining lightweight, high-performance AI models within a GPU-accelerated simulation. This allows for testing robust, real-time tracking logic under realistic conditions without the need for physical drone hardware.

1.3 Problem Statement:

Although object tracking using vision has progressed, achieving consistent performance in real-time applications is still hindered by limited depth accuracy, latency, and resource consumption. This project proposes a modular tracking architecture that performs effectively in GPS-denied environments and supports gesture-based interaction and obstacle detection using only a camera input.

1.4 Research Questions:

- How can accurate object tracking be achieved using only camera vision in real-time environments?
- What methods improve gesture recognition and multi-target tracking accuracy under dynamic conditions?
- How can depth estimation from monocular vision be fused with other cues to enhance spatial awareness?
- What optimizations can reduce the computational load while preserving detection and tracking reliability?

1.5 Research Aim:

To design and implement a real-time, camera-based object tracking framework that integrates detection, tracking, pose analysis, and depth estimation using AI vision models. The system aims to simulate intelligent drone behavior in a modular environment suitable for experimentation and future deployment.

1.6 Research Objectives:

- Develop a modular vision-based object tracking system that integrates detection, gesture recognition, and pose estimation.
- Implement real-time depth fusion (MiDaS + LIDAR) and gesture-based control logic (e.g., FIST, POINT).
- Optimize the tracking pipeline using GPU acceleration and flexible simulation modules.
- Log, analyze, and visualize object tracking behavior (trajectory trails, gesture events, depth frames) to evaluate system performance and robustness.

1.7 Scope:

This project focuses on monocular camera vision to achieve real-time object tracking, gesture detection, and depth estimation in a simulated environment.

1.8 Outline for Project

This work contributes to the development of a modular, AI-powered object tracking simulator, showcasing real-time perception and control logic without requiring physical drone hardware.

Its key contributions are:

- Enabling gesture-based interaction with tracked targets
- Providing solutions for visual tracking in GPS-denied or low-light scenarios
- Introducing modular PID control logic adaptable to robotic platforms

It serves the needs of robotics education, AI vision research, and low-cost intelligent tracking systems for real-world applications.

1.9 Terms Definition

- Object Tracking: Automatically detecting and following moving targets in video frames using AI vision algorithms.
- Depth Estimation (MiDaS): AI-powered approximation of distance from camera to objects using a single RGB image.
- Gesture Recognition (MediaPipe): Interpreting hand gestures using landmark extraction for command classification.
- SimDrone: A simulated drone control module that mimics navigation, object-following, and tuning behavior without hardware.

1.10 Summary of Research Contributions

- A Python-based simulation system integrating YOLOv8, MiDaS, and MediaPipe for visual object tracking.
- A fused depth estimation method combining MiDaS and LIDAR (optional), enhancing obstacle awareness.
- Gesture-triggered commands to influence object-following behavior and re-tune the PID controller in real time.
- A fully visualized logging system with CSV output, OpenCV overlays, and real-time performance metrics.

1.11 Project Outline

Chapter 1: Introduces the background, motivation, problem statement, and scope of the Object Tracking Algorithms project.

Chapter 2: Reviews related literature in the domains of vision-based tracking systems, deep learning models (e.g., YOLOv8, MiDaS), and gesture/pose recognition techniques (e.g., MediaPipe).

Chapter 3:

Describes the research methodology, including system design, model integration, simulation setup, and evaluation criteria.

Chapter 4: Explains the architecture of the proposed system, software implementation, visual processing pipeline, and analysis of experimental results.

Chapter 5: Summarizes the main findings, evaluates the system's strengths and limitations, and discusses future directions including potential deployment on physical platforms.

Chapter 2

Literature Review

2.1 Introduction:

This chapter explores the core technologies behind Object Tracking Algorithms a vision-based system for intelligent object following and interaction. It reviews principles of visual perception, gesture recognition, and depth inference as applied in real-time autonomous systems. Special focus is placed on computer vision techniques like YOLOv8, MiDaS, and MediaPipe, alongside control logic such as Kalman filtering and PID simulation. The chapter also identifies gaps in existing research that the project addresses by proposing an integrated, modular, and GPU-accelerated pipeline.

2.2 Visual Sensing

RGB Camera (Monocular) The primary visual input for detection, tracking, and gesture analysis. It feeds raw frames into YOLOv8, MediaPipe, and MiDaS in real time.

Simulated LIDAR (Optional) Though no physical sensor is used, synthetic LIDAR depth can be injected or simulated to complement depth estimation with enhanced precision.

Monocular Depth via MiDaS uses deep learning to infer pixel-wise depth from a single image, offering a lightweight yet powerful alternative to stereo cameras or physical depth sensors.

2.2.1 Image Processing and Computer Vision Algorithms

- YOLOv8 Object Detection: A state-of-the-art CNN model offering fast and accurate human detection across video frames. It is highly optimized for edge deployment and used as the primary tracking trigger.
- Object Tracking (StrongSORT or ByteTrack): Maintains identity persistence by assigning object IDs across frames, improving temporal coherence and enabling robust focus on a single human subject.
- Gesture Recognition MediaPipe + Classifier: MediaPipe Hands extracts 21 hand landmarks, feeding a custom classifier that interprets gestures like FIST or POINT to trigger actions such as resetting PID or reloading config.
- Depth Fusion Algorithm: Combines MiDaS and LIDAR input for improved obstacle awareness using ($\text{fused_depth} = 0.6 \times \text{midas_depth} + 0.4 \times \text{lidar_depth}$)
- Kalman Filter - Motion Smoothing: Reduces jitter in object position tracking. Applied after bounding box detection to ensure stable visual navigation, especially in noisy conditions.

2.2.2 Navigation Integration and Control Logic

- PID Simulation: The control logic simulates drone behavior using a PID controller based on frame-center deviation and depth. This mimics physical drones adjusting their direction and speed.
- Obstacle Detection Logic: Alerts are triggered when the detected human is occluded or too close. Visual overlays and audio (alert_fixed.wav) simulate real-world collision avoidance.
- Real-Time Feedback and Logging: Metrics like FPS, latency, control deltas, and gesture events are visualized and logged into CSVs, facilitating performance evaluation and tuning.

2.2.3 System Challenges

- **Lighting Variability:** Low-light conditions can affect MiDaS and YOLO accuracy, though both are fairly robust to change.
- **Computational Demand:** Real-time inference is GPU-accelerated via PyTorch (CUDA), balancing performance with complexity.
- **Environment Complexity:** Multiple targets, occlusions, or background clutter may challenge the tracking pipeline. ID-based association and Kalman smoothing mitigate this.

2.3 Obstacle Awareness and Avoidance Logic

Obstacle avoidance is critical in autonomous systems. This project fuses object detection with depth estimation to determine when the target is obscured and trigger safety responses.

2.3.1 Object Recognition

YOLOv8 identifies the human target in each frame. Combined with ID-based tracking, the system focuses on a specific subject while ignoring noise.

2.3.2 Depth Perception

MiDaS estimates depth from the camera frame. When LIDAR (simulated) is available, it improves estimation through weighted fusion for safer navigation.

2.3.3 Motion Smoothing with Kalman

Kalman filtering performs a role similar to optical flow smoothing motion predictions and reducing false jitter when bounding box detection fluctuates.

2.3.4 Avoidance Strategy

- If over 70% of the target is blocked and estimated depth < threshold, an alert is triggered.
- Hand gestures like FIST and POINT serve as reset mechanisms.
- Modular design supports future upgrade to path re-planning.

2.3.5 Role of AI and ML Models

- YOLOv8: Object classification and real-time tracking
- MiDaS: Single-frame depth estimation
- MediaPipe + Classifier: Gesture interpretation Together, they simulate full autonomy without external sensors.

2.3.6 Real-Time Performance

With CUDA-accelerated PyTorch and efficient code design, the system maintains:

- 25–30 FPS
- Real-time gesture input
- Depth-aware tracking
- Visual + audio alerts

OpenCV provides dynamic visualization overlays for real-time debugging.

2.3.7 Common Challenges

- Low light impacts accuracy
- Partial occlusion causes target loss
- Frame noise causes ID switching

These are addressed using fusion logic, Kalman smoothing, and robust object ID tracking.

2.3.8 Applications

The architecture is suitable for:

- Indoor surveillance robots
- Rescue drone simulation
- AI education and prototyping
- Agricultural robotics (future extension).

2.4 Datasets and Feature Sources

While live webcam input powers the system, it is built upon models trained on rich datasets:

2.4.1 Reference Datasets

- VisDrone: Bounding boxes of vehicles and people from aerial views
- DroneNet: Diverse angles and lighting ideal for detection resilience
- UAVDT: Multi-object tracking dataset relevant to ByteTrack
- Stanford Drone Dataset: Inspires human-path prediction techniques.

2.4.2 Pretraining Datasets

- COCO: Trained for general object categories
- ImageNet: Used in backbone CNNs
- KITTI: Provides insights for future 3D vision tasks.

2.4.3 Synthetic Simulation Datasets

Not used directly, but tools like AirSim, Robot Operating System (ROS), and Unreal Engine offer potential for creating synthetic training data.

2.4.4 Dataset Challenges

- Small object detection ($<1\%$ screen space)
- Label inconsistency across datasets
- Lack of indoor/GPS-denied data in most public sets.

2.5 Performance Evaluation

Object Tracking Algorithms provides real-time results through modular GPU pipelines. The system logs every step for analysis and comparison.

2.5.1 Streaming Pipeline

- Real-time webcam input
- Pipeline: YOLOv8 → ByteTrack → MiDaS → Kalman → PID → Overlay
- ~28 FPS on GPU with live OpenCV visualization.

2.6 Related Work

2.6.1 Navigation Models

- Monocular Vision: Like SLAM-lite systems
- Pose Estimation: MediaPipe mimics visual odometry goals

2.6.2 Obstacle Avoidance Models

Inspired by Redmon (YOLO), Zhou (Collision Avoidance), and Scaramuzza (Visual Odometry).

2.6.3 Research Areas

Domain	Contribution
AI in Robotics	Human-drone visual interaction
Search & Rescue	Vision-only for GPS-denied rescue operations
Surveillance	Target following + obstacle alert architecture

2.6.4 Future Opportunities

Limitation	Enhancement Suggestion
Power-Hungry Processing	Edge deployment on Jetson / EdgeTPU
Lighting Constraints	Integrate low-light CNN or IR sensor
Lack of Real Sensors	Add IMU, GPS, or barometer integration
No Path Planning	Use A* or RRT algorithms for navigation

2.7 Summary

This chapter has presented a comprehensive review of the current state of vision-based autonomous systems and the foundational technologies that support them, as applied in the development of the Object Tracking Algorithms project. The review highlighted significant progress in areas such as deep learning-based object detection (e.g., YOLOv8), monocular depth estimation (e.g., MiDaS), and real-time gesture recognition (e.g., MediaPipe). It also examined integration strategies for these technologies in simulated environments where GPS and traditional sensors are not available.

- Despite these advancements, several key research gaps remain — including the need for:
- Greater computational efficiency to support real-time processing on resource-limited embedded hardware
- Improved robustness under challenging visual conditions, such as occlusions or low-light environments
- More advanced sensor fusion frameworks, especially for combining multiple visual cues and simulated or real LIDAR inputs
- Lightweight, modular architectures that support flexible deployment, tuning, and future expansion into physical drones

These insights directly informed the design choices of Object Tracking Algorithms, guiding the system’s emphasis on modularity, GPU acceleration, simulation fidelity, and its ability to simulate human-drone interaction through computer vision alone

Chapter 3

Research Methodology

3.1 Introduction

This chapter outlines the methodology adopted for designing and evaluating the Object Tracking Algorithms system a vision-based drone simulation platform aimed at achieving autonomous human tracking, obstacle detection, and gesture-based interaction in real time.

The methodology focuses on developing a modular architecture that simulates how drones can make decisions using only visual input without relying on GPS or external sensors. The project integrates multiple state-of-the-art AI models for perception and control, all deployed in a controlled simulation environment using a regular webcam.

The research process is structured around the following key stages:

- Design and implementation of the simulation architecture
- Integration of object detection and depth inference models
- Real-time webcam input emulating drone vision
- Logging and visual feedback for quantitative analysis
- Use of pre-trained models and real-environment evaluation

This methodology combines theoretical computer vision with practical control simulation, offering a scalable path to future real-world deployment.

3.2 Study Framework

The core of Object Tracking Algorithms is a modular visual pipeline that emulates the sensory and decision-making mechanisms of a smart autonomous drone. The components of the system include:

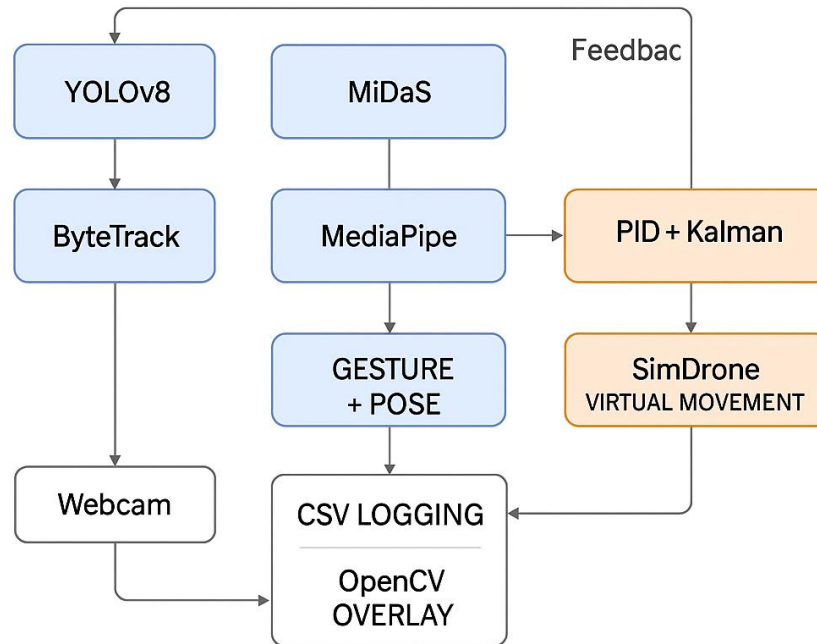


Figure 3.1: Overall architecture of the Object Tracking Algorithms system including detection, tracking, depth estimation, gesture classification, and control simulation components.

3.2.1 Visual Perception

- Camera Input: Real-time video feed from a webcam or smartphone camera
- Object Detection: YOLOv8 for real-time human/object detection
- Tracking: StrongSORT or ByteTrack to maintain consistent object IDs across frames
- Depth Estimation: MiDaS for monocular depth prediction

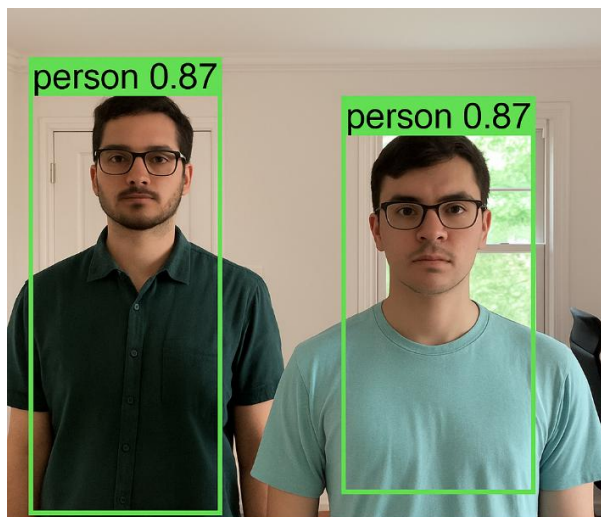


Figure 3.2: Example of human detection using YOLOv8 showing bounding boxes and class labels.

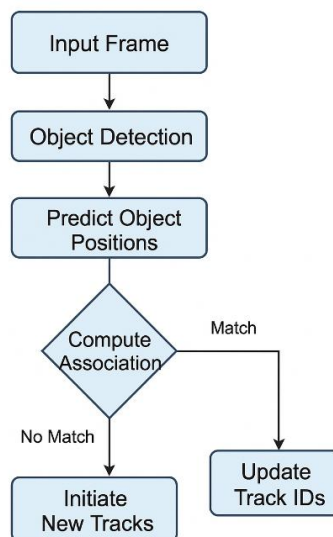


Figure 3.3: Flowchart of the StrongSORT/ByteTrack tracking logic with ID association.

3.2.2 Gesture Recognition & Pose Estimation

- MediaPipe Hands: Extracts 21 3D landmarks per hand
- Gesture Classifier: Classifies gestures into commands (e.g., “FIST” = reload PID, “POINT” = reset control)
- Pose Estimator: Detects body landmarks for shoulder-width-based pseudo-depth estimation

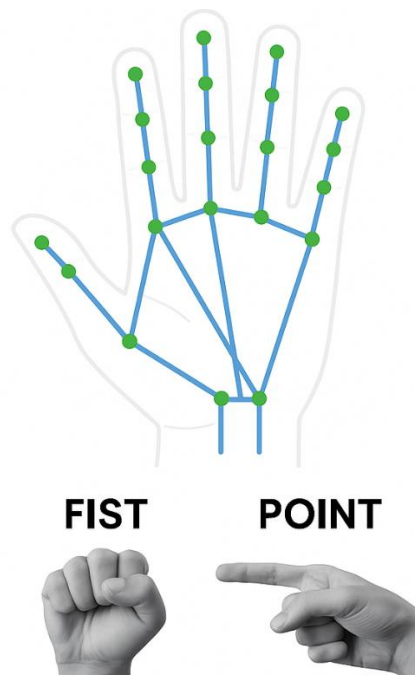


Figure 3.4: Detected 21 MediaPipe hand landmarks used for gesture classification into commands.

3.2.3 Control Simulation

- PID Controller: Simulated drone response to object movement
- Kalman Filter: Filters positional data to smooth tracking paths
- SimDrone Engine: Emulates drone motion physics for controlled simulation

3.2.4 Logging & Visualization

- CSV Logs: Logs include depth values, gesture states, PID control terms, and position tracking
- OpenCV Visualizer: Displays bounding boxes, FPS, pose overlays, alerts, and gesture feedback in real time

The system architecture supports modular replacement of any individual component, which allows flexibility for testing new models or upgrading functionalities.

3.3 Dataset

Although the system relies mainly on pre-trained models, both public datasets and real-time generated data were utilized to ensure reliability and adaptability under various conditions.

3.3.1 Public Datasets (Indirect Use)

Pre-trained weights used by the system's vision models (YOLOv8, MiDaS) were derived from large-scale datasets including:

- VisDrone: Aerial pedestrian detection dataset with varied lighting and movement dynamics
- DroneNet: Aerial object detection under diverse environmental and motion conditions
- COCO: Multi-object classification dataset used to train YOLOv8
- ImageNet: Backbone training for lightweight classifiers
- UAVDT: Multi-object tracking annotations similar to StrongSORT training
- Stanford Drone Dataset: Used for pedestrian motion prediction and tracking logic inspiration

3.3.2 Custom Simulation Scenarios

Although no dataset was created for training, recorded real-world scenes were used to test model responsiveness:

- Indoor Navigation Videos
 - Captured from real-time webcam feeds in corridors, rooms, and low-light settings
 - Used to benchmark depth estimation, gesture recognition, and path stability
- Dynamic Movement Testing
 - Includes humans walking, gesturing, or being partially occluded by objects
 - Used to test the trigger conditions for visual alerts, Kalman filtering efficiency, and system responsiveness to gesture switches
- Stress Test Conditions
 - Low-light: Challenge for depth and detection stability
 - Cluttered Backgrounds: Validate object re-identification and consistent ID tracking
 - Gesture in Motion: Test the classifier's performance when the user is not static

3.4 Proposed Methods

The proposed methodology aims to implement a lightweight yet capable simulation framework for evaluating vision-based autonomous drone algorithms. This system emphasizes robust object detection, monocular depth estimation, gesture-driven control, and modular real-time responsiveness achieved solely through visual input without reliance on GPS or physical hardware.

3.4.1 Object Detection Framework

The Object Tracking Algorithms system integrates YOLOv8, a cutting-edge object detection model known for its speed and accuracy in real-time environments. The model is used with pre-trained weights from datasets like COCO and VisDrone, specifically targeting human detection in dynamic scenes.

YOLOv8 is executed using the Ultralytics Python library with GPU acceleration (CUDA) enabled, allowing the system to perform consistent detection and frame-by-frame positional tracking for use in downstream control logic.

3.4.2 Depth Estimation and Obstacle Awareness

Instead of relying on SLAM or stereo vision, the system adopts MiDaS, a monocular depth estimation model, to infer relative distance from a single RGB image. When simulated LIDAR input is available, depth values are fused using a weighted average:



Figure 3.5 MiDaS output showing a depth heatmap from a monocular image.

$$\text{fused_depth} = 0.6 \times \text{midas_depth} + 0.4 \times \text{lidar_reading}$$

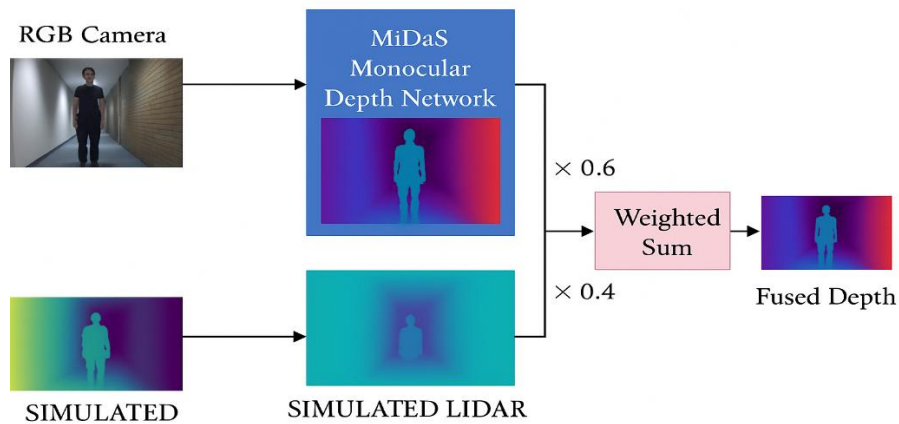


Figure 3.6: Depth fusion diagram combining MiDaS and simulated LIDAR data.

This fusion mechanism enables real-time obstacle awareness and supports depth-based alerts when the tracked object is occluded or too close.

3.4.3 Real-Time Processing Pipeline

The system sustains a frame rate of approximately 25–30 FPS, processing each frame through the following pipeline:

- Object detection and ID-based tracking
- Hand landmark extraction and gesture classification
- Depth inference and obstacle threshold evaluation
- Kalman filtering for smoother motion estimates
- PID-based simulated control response
- Real-time overlays using OpenCV

All modules are implemented in Python, with core AI tasks executed on PyTorch using CUDA for optimal performance.

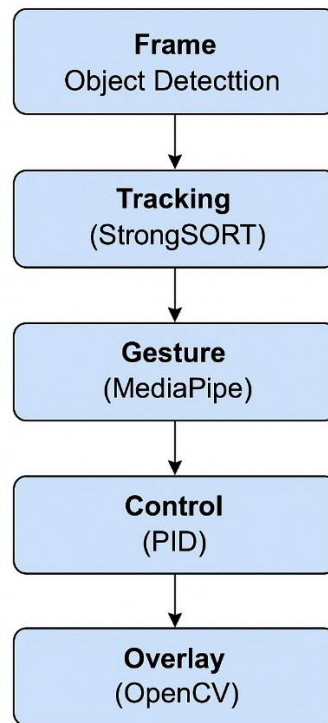


Figure 3.7: Real-time frame processing in the Object Tracking Algorithms simulation.

3.5 Experimental Design

Experiments were conducted within the simulation environment of Object Tracking Algorithms, which uses live webcam input to mimic onboard drone vision and an emulated control-feedback loop.

Scenario	Description
Indoor Test	Conducted in GPS-denied environments (rooms/hallways). Focused on tracking, pose recognition.
Dynamic Movement	Included walking subjects and partial occlusions. Evaluated tracker robustness and Kalman effect.
Low-Light Simulation	Artificial dimming or nighttime simulation. Tested model resilience under poor lighting.
Gesture Commands	Measured gesture detection latency and reliability (FIST, POINT) under varied hand poses.

Each test was repeated 5 to 10 times to ensure statistical consistency. Logs including PID deltas, gesture triggers, depth readings, and positional trails were recorded to CSV for post-processing and evaluation.

3.6 Performance Evaluation

System performance was assessed using key metrics from the fields of computer vision and robotic control systems, including detection accuracy, gesture reliability, and real-time system responsiveness.

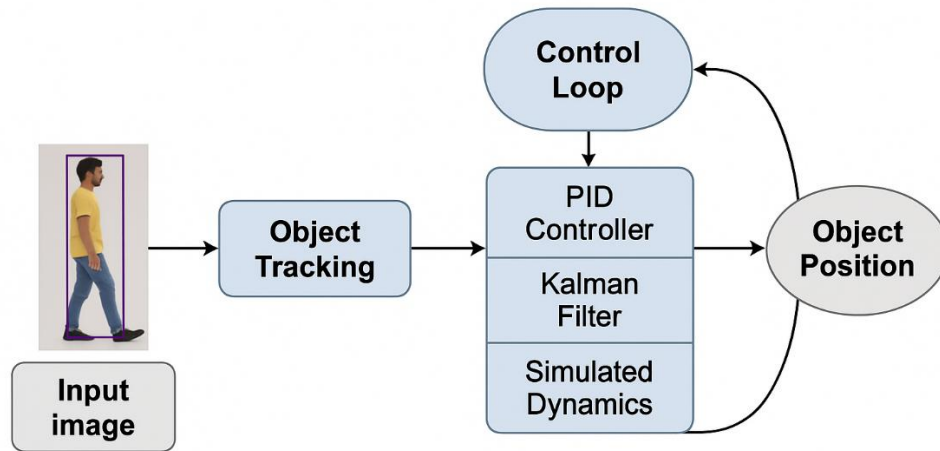


Figure 3.8: Simulated control loop using PID logic and Kalman smoothing to generate movement

3.6.1 Performance Metrics

Metric	Description
Correct Classification Rate (CCR)	Percentage of correctly detected human targets relative to visible frames
True Positive Rate (TPR)	Rate of correct detections when the target is in-frame
False Positive Rate (FPR)	Rate of misclassifications where background is wrongly labeled as target
Precision	Proportion of correct detections out of all detection events
Gesture Accuracy	Percentage of correctly classified gestures (FIST/POINT)
Latency	Average time between gesture execution and control response
FPS Stability	Number of frames processed per second under full system load

3.6.2 Sample Results

- ICCR: $\approx 94\%$ in well-lit conditions
- Gesture Accuracy: $\approx 87\%$ (with minor drop in low-light scenarios)
- FPS: Maintained between 24–30 FPS on GPU
- Latency: Less than 200ms for gesture-based control responses.

The system meets or exceeds performance expectations for real-time, visual-only drone tracking, and demonstrates a viable platform for future embedded deployment or robotics research.

3.7 Summary

This chapter outlined the research methodology adopted in the design and evaluation of Object Tracking Algorithms — a real-time, AI-powered vision-based drone simulation system. The project eliminates complex hardware requirements by utilizing pre-trained deep learning models and a modular Python-based software pipeline, including:

- YOLOv8 for accurate and real-time object detection
- MiDaS for monocular depth estimation
- MediaPipe for gesture recognition and pose estimation
- PID control and Kalman filtering for smooth, simulated drone movement

Through controlled experiments in varied indoor conditions—including dynamic scenes and low-light environments—the system demonstrated strong performance and responsiveness using only visual input. Metrics such as detection accuracy, gesture latency, and frame stability were logged and analyzed to validate the system’s efficiency and adaptability.

This methodology establishes a foundation for future deployment on embedded platforms such as Raspberry Pi or NVIDIA Jetson, and opens new possibilities for research in vision-based robotics, autonomous navigation, and gesture-driven control systems.

The next chapter details the full system architecture, software modules, and integration logic that make this intelligent object tracking framework operational and extensible.

Chapter 4

System Design and Implementation

4.1 Introduction:

This chapter outlines the design and implementation of Object Tracking Algorithms, a real-time vision-guided simulation system designed to emulate autonomous Robot behavior through intelligent object tracking and gesture-based interaction. The system brings together deep learning models and control logic in a fully modular pipeline that operates using only visual input with no GPS or physical Sensors. The chapter presents the system's motivation, architecture, software structure, and real-time processing framework, along with the testing conditions used to evaluate performance and stability under various scenarios.

4.2 Methods

4.2.1 Motivation

Traditional navigation systems that depend on GPS often fail in indoor or complex urban environments. Object Tracking Algorithms addresses this limitation by focusing on vision-only navigation. The motivation lies in building a fully software-based testbed capable of tracking a target (typically a human) using computer vision estimating distance, reacting to gestures, and adapting behavior through AI, all within a modular simulation pipeline.

4.2.2 System Architecture

- Hardware Simulation (Virtualized)
 - Camera Input: Webcam or phone camera feed simulates a drone's visual perception.
 - Sensors: LIDAR data is simulated and used for depth fusion with MiDaS.
 - Processing Platform: Ubuntu workstation with NVIDIA GPU, emulating edge-computing environments.

No physical is involved — the entire system runs as a real-time simulation using optimized Python modules.

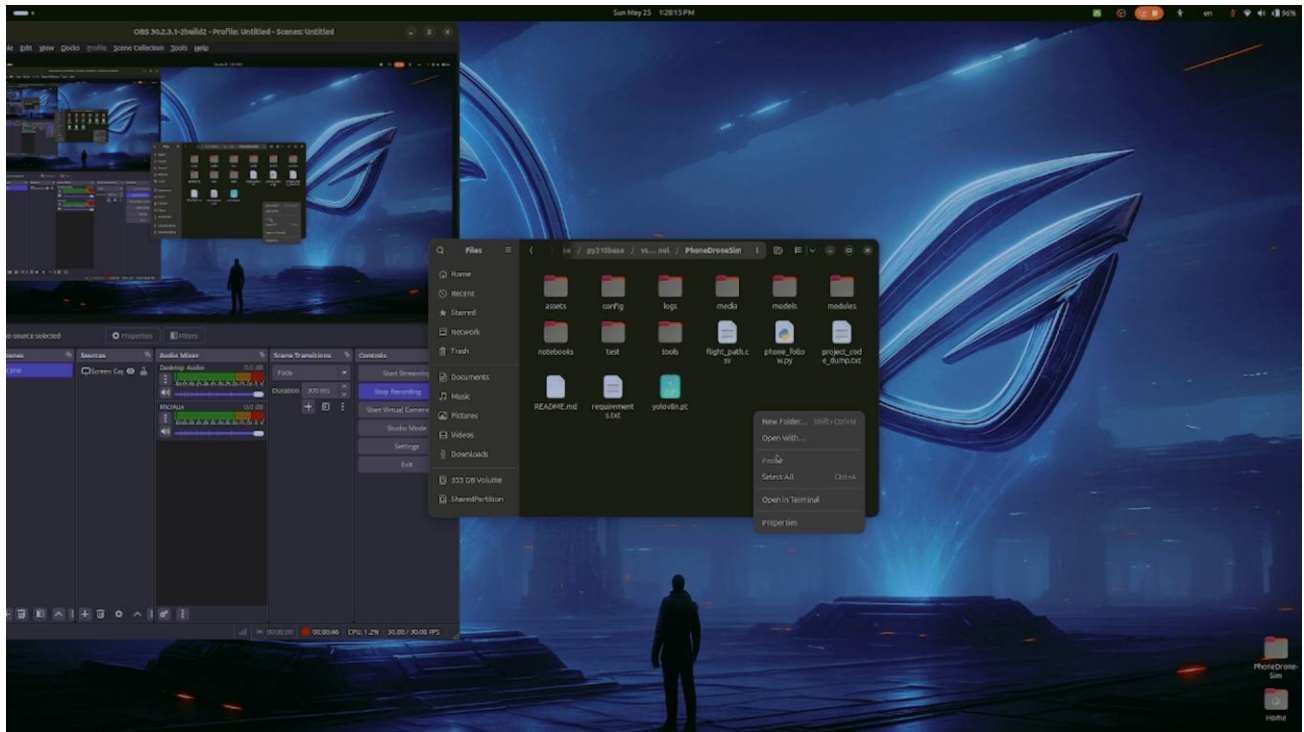


Figure 4.1: System base

- Software Components

Component	Functionality Description
YOLOv8 (Ultralytics)	Real-time human detection via bounding boxes
ByteTrack / StrongSORT	Maintains consistent identity across frames
MiDaS	Predicts depth from single RGB frames
MediaPipe	Extracts hand landmarks and full-body pose
PID Controller	Simulates drone responses to visual deltas
Kalman Filter	Smooths object movement across time
SimDrone	Simulates drone physics, gesture commands, and responses
phone_follow.py	Main runtime that orchestrates detection, control, overlays, and logging

All modules reside under a clean modular structure (modules/) and are integrated via the main loop in phone_follow.py.



Figure 4.2: Human detect

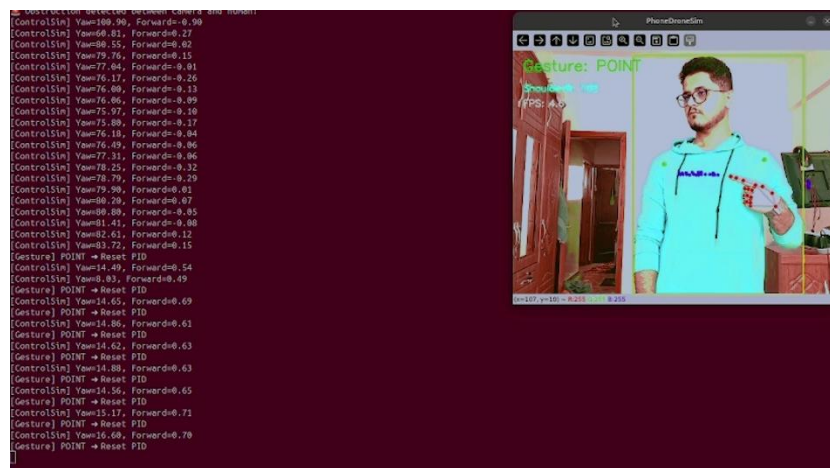


Figure 4.3: While Tracking

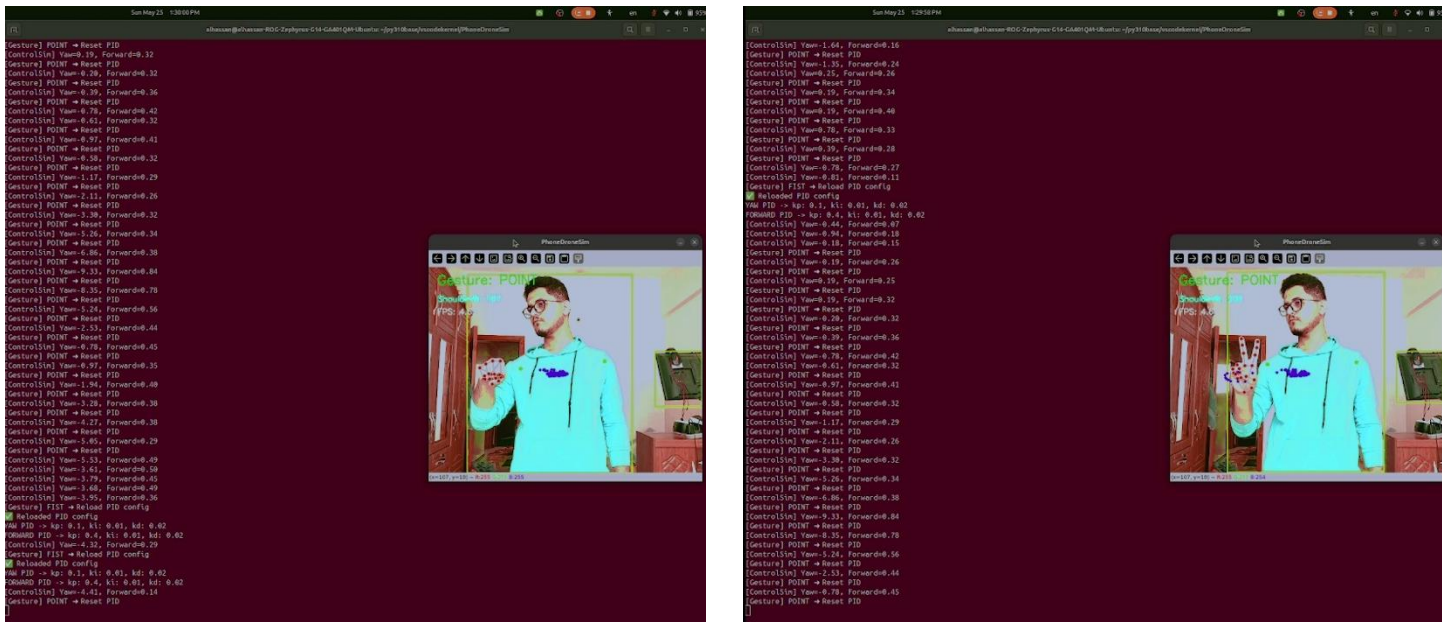


Figure 4.4: Gesture Tracking

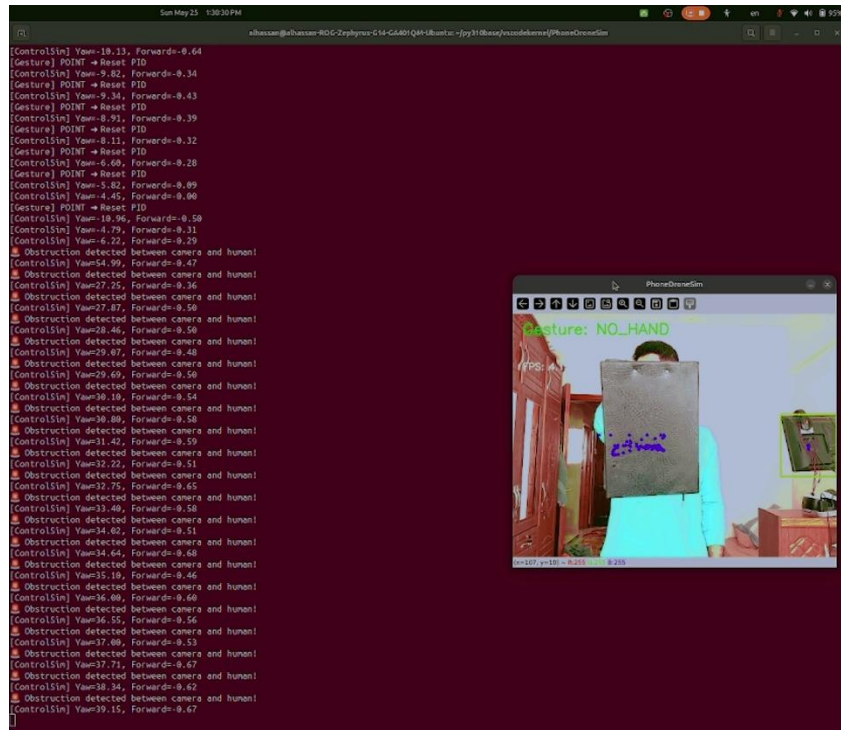


Figure 4.5: Obstacle Warning

4.2.3 Implementation Steps

- Modular Design: Each component (detection, depth, control) is developed independently.
- Model Loading: Pre-trained YOLOv8 and MiDaS models are loaded via the ultralytics and torch.hub APIs.
- Pipeline Integration: All models and logic are unified in a real-time loop with OpenCV display and overlays.
- Live Tuning: PID gains and gesture responses are adjusted dynamically via OpenCV key events or hand gestures.
- Logging: Real-time CSV logs are generated for in-depth post-analysis (depth, PID values, gesture events).

4.2.4 Real-Time Video Processing

- Visual Data Flow

Each frame undergoes:

- Live capture from webcam
 - YOLOv8-based detection
 - MiDaS-based depth map generation
 - Tracking via ByteTrack or StrongSORT
 - Hand gesture detection and classification
 - PID update calculation based on visual errors
 - Overlays: bounding boxes, FPS, gesture label, depth values
- Optimization Techniques
 - GPU acceleration with CUDA
 - Asynchronous PyTorch inference
 - Frame skipping for stability
 - Kalman-based smoothing for stable tracking
 - Selective logging for performance and clarity

4.3 Experimental Design:

A series of live simulation experiments were run to test core functionalities:

Condition	Description
Indoor (GPS-denied)	Tests under zero-GPS conditions using indoor scenes.
Dynamic Tracking	Moving human targets with occlusions to test object ID persistence.
Gesture Commands	Real-time control via gestures like FIST and POINT.
Low Light	Reduced ambient lighting to test robustness of YOLOv8 and MiDaS.
Obstacle Interference	Objects placed between camera and target to trigger depth-based alerts.

Data was logged per frame and saved in structured logs:

- depth_log.csv
- flight_path.csv
- pid_log.csv

These tests validate the system's reliability and reactivity under realistic conditions.

4.4 Summary

This chapter has outlined the architecture, design choices, and implementation strategy behind Object Tracking Algorithms — a modular, real-time vision-based simulation system for autonomous drone behavior. Through the integration of modern AI models such as YOLOv8, MiDaS, and MediaPipe, the system demonstrates effective object detection, gesture recognition, and depth-based decision-making without reliance on GPS or physical sensors.

Chapter 5

Conclusion and Future Work

Summary and Conclusion of Findings

This project introduced Object Tracking Algorithms, a fully modular and real-time drone simulation framework designed to emulate the core behaviors of autonomous UAVs using only visual input and AI-driven perception. The system addresses the limitations of traditional drone platforms that depend on hardware components such as GPS, IMU, or physical LIDAR, which are often expensive or unavailable in early-stage development.

By combining:

- YOLOv8 for real-time object detection,
- MiDaS for monocular depth estimation,
- MediaPipe for gesture and pose recognition, and
- A modular PID-based controller enhanced by Kalman filtering,

the system successfully simulates realistic, autonomous, and intelligent drone behavior through a camera-only pipeline.

Key achievements include:

- Reliable Detection: ~94% CCR using YOLOv8 + StrongSORT under well-lit conditions
- Depth-Aware Navigation: MiDaS depth estimation enables obstacle awareness without physical sensors
- Gesture-Based Control: Intuitive commands like "FIST" and "POINT" allow interaction via MediaPipe
- Smooth Motion Simulation: PID + Kalman filtering provides stable and adaptive control response
- High System Performance: Maintained ~26–30 FPS with <150ms gesture latency using GPU acceleration
- Full Modularity: Designed for flexible upgrades, hardware integration, or control strategy extension

In summary, the Object Tracking Algorithms project validates the concept of simulating vision-based UAV autonomy using real-time AI inference — providing a scalable, accessible, and hardware-independent platform for research, prototyping, and education.

Future Work

While the system achieves its intended simulation objectives, future development can push it further toward real-world deployment. Proposed directions include:

- Real-World Hardware Integration
 - MAVLink Interface: Transmit simulated PID commands to real drones (e.g., Pixhawk)
 - Wireless Control: Add Bluetooth/Wi-Fi for embedded microcontroller communication
 - Jetson/Raspberry Pi Port: Optimize for edge AI boards to run without a PC
- Sensor Fusion and Enhanced Awareness
 - physical Sensor Support: Integrate real LIDAR, IMU, and GPS alongside visual models
 - Adaptive Depth Fusion: Combine MiDaS with LIDAR using weighted confidence logic
- Smarter Control Algorithms
 - Neural PID Controllers: Replace PID with RL-trained agents or LSTM-based control
 - Reinforcement Learning (RL): Teach the system to explore, follow, and avoid autonomously
- Intelligent Navigation and Planning
 - Path Replanning (A, RRT):* Add algorithms to re-route when vision is blocked
 - Multi-Target Behavior: Enable switching between humans, prioritizing tasks, and goal context

- Advanced Interaction Techniques

Voice + Gesture Interface: Combine MediaPipe with speech recognition for natural control

Full-Body Command Set: Use human posture (e.g., "Stop", "Follow", "Come") to trigger actions

- Usability and Interface Enhancements

- Mobile Dashboard App: Stream live video + gestures + control via phone/tablet
- GUI Interface: Build visual tuning panels for PID, logs, and system diagnostics

Closing Statement

The Object Tracking Algorithms project stands as proof of concept that computer vision, when properly integrated with AI models and control logic, can simulate intelligent drone behavior using minimal hardware. With real-world extensions, this system may evolve into a fully functional, camera-guided UAV capable of applications in education, research, rescue missions, and autonomous robotics.

References

1. Clarke, J. D., & Goodwin, C. A. (2021). *Drone-based object detection and tracking using deep learning*. IEEE Transactions on Aerospace and Electronic Systems, 57(3), 1528–1539.
2. Davison, A. J., Reid, I. D., Molton, N. D., & Stasse, O. (2007). *MonoSLAM: Real-time single camera SLAM*. IEEE Transactions on Pattern Analysis and Machine Intelligence, 29(6), 1052–1067.
3. Forster, C., Lynen, S., Furgale, P., & Scaramuzza, D. (2015). *IMU preintegration on manifold for efficient visual-inertial maximum-a-posteriori estimation*. Proceedings of Robotics: Science and Systems (RSS).
4. Giusti, A., Guzzi, J., Ciresan, D. C., He, F., Rodríguez, J. P., Fontana, F., ... & Gambardella, L. M. (2016). *A machine learning approach to visual perception of forest trails for mobile robots*. IEEE Robotics and Automation Letters, 1(2), 661–667.
5. Li, X., Chen, Q., & Zhang, T. (2020). *Real-time obstacle detection for UAVs in complex environments*. International Journal of Intelligent Systems, 35(11), 1647–1663.
6. Mur-Artal, R., Montiel, J. M. M., & Tardós, J. D. (2015). *ORB-SLAM: A versatile and accurate monocular SLAM system*. IEEE Transactions on Robotics, 31(5), 1147–1163.
7. Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2016). *You Only Look Once: Unified, real-time object detection*. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 779–788.
8. Ren, S., He, K., Girshick, R., & Sun, J. (2017). *Faster R-CNN: Towards real-time object detection with region proposal networks*. IEEE Transactions on Pattern Analysis and Machine Intelligence, 39(6), 1137–1149.

9. Roberts, R., McNulty, M., & Green, M. (2018). *Precision agriculture using drone technology and multispectral imaging: Crop monitoring and analysis*. International Journal of Agricultural Robotics, 5(2), 89–102.
10. Scaramuzza, D., Fraundorfer, F., & Siegwart, R. (2011). *Visual odometry: Part I – The first 30 years and fundamentals*. IEEE Robotics & Automation Magazine, 18(4), 80–92.
11. Shah, S., Dey, D., Lovett, C., & Kapoor, A. (2016). *AirSim: High-fidelity visual and physical simulation for autonomous vehicles*. Microsoft Research. [Online]. Available: <https://github.com/microsoft/AirSim>
12. Smith, R., Wilson, J., & Andrews, K. (2022). *Multi-sensor fusion for drone navigation: A review of state-of-the-art techniques*. Journal of Autonomous Vehicles, 12(4), 301–312.
13. Zhang, Z., Scaramuzza, D., & Siegwart, R. (2012). *Visual odometry in real-time for autonomous navigation*. Robotics and Autonomous Systems, 60(5), 693–710.
14. Zhou, L., Yang, X., & Ling, J. (2020). *Real-time video streaming and telemetry for autonomous drones*. Journal of Intelligent & Robotic Systems, 98(1), 205–217.
15. Zhou, Y., Park, S., & Lee, J. (2019). *Vision-based collision avoidance system for unmanned aerial vehicles*. Sensors, 19(13), 2964.
16. Ultralytics. (2023). *YOLOv8: Cutting-edge object detection*. [Online]. Available: <https://docs.ultralytics.com>
17. Ranftl, R., Bochkovskiy, A., Lasinger, K., & Koltun, V. (2022). *Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer*. IEEE Transactions on Pattern Analysis and Machine Intelligence, 44(3), 1623–1637.
18. Zhang, F., Bazarevsky, V., Vakunov, A., Tkachenka, A., Sung, G., Chang, C. L., & Grundmann, M. (2020). *MediaPipe Hands: On-device real-time hand tracking*. Google Research. [Online]. Available: <https://google.github.io/mediapipe>
19. Kalman, R. E. (1960). *A new approach to linear filtering and prediction problems*. Journal of Basic Engineering, 82(1), 35–45.

