



PEOPLES'S DEMOCRATIC REPUBLIC OF ALGERIA
MINISTRY OF HIGHER EDUCATION AND SCIENTIFIC RESEARCH

Blida 01 University
Institute of aeronautics and spatial studies
Département Construction



Manuscript In partial fulfilment of the requirements for the Degree of Master



Specialty: Avionics

GESTURE CONTROL OF A QUADROTOR UAV

By



Meriam BOUABBA

Kh. Choutri
M. Lagha

Doctor, IAES. De Blida
Doctor, IAES. De Blida

Advisor
Vice-Advisor

Blida, juillet 2025

Abstract:

Natural interaction systems, particularly those based on computer vision, provide intuitive and effective communication between humans and drones. Gesture control allows for direct piloting of the drone, while person tracking ensures that the drone maintains visual contact with the operator. One function without the other is incomplete: gesture control is useless if the drone cannot see the user, and simple tracking offers no means of command. This project aims to combine these two functionalities by developing an intelligent system where a gesture recognition model and a person detection model work in tandem. The objective is to create a seamless, robust, and non-restrictive human-drone interface, where the user's body becomes the primary control device.

Key words:

Human-drone interaction, gesture recognition, person tracking, computer vision, drone control, natural interaction.

Resume:

Les systèmes d'interaction naturelle, notamment ceux basés sur la vision par ordinateur, offrent une communication intuitive et efficace entre l'humain et le drone. Le contrôle par gestes permet de piloter le drone de manière directe, tandis que le suivi de personne (tracking) assure que le drone maintient un contact visuel avec l'opérateur. L'un sans l'autre est incomplet : un contrôle gestuel est inutile si le drone ne peut pas voir l'utilisateur, et un simple suivi n'offre pas de contrôle. Ce projet vise à combiner ces deux fonctionnalités en développant un système intelligent où un modèle de reconnaissance de gestes et un modèle de détection de personne fonctionnent de concert. L'objectif est de créer une interface homme-drone fluide, robuste et non restrictive, où le corps de l'utilisateur devient le véritable dispositif de contrôle.

Mots-clés :

Interaction homme-drone, reconnaissance de gestes, suivi de personne, vision par ordinateur, contrôle de drone, interaction naturelle.

ملخص:

تتيح أنظمة التفاعل الطبيعي، خاصة تلك المعتمدة على الرؤية الحاسوبية، تواسلاً سهلاً وفعالاً بين الإنسان والطائرة بدون طيار. يسمح التحكم عبر الإيماءات بتوجيه الطائرة بشكل مباشر، بينما يضمن تتبع الأشخاص الحفاظ على اتصال بصري مستمر مع المشغل. إحدى الوظيفتين لا تكتمل بدون الأخرى: فالتحكم بالإيماءات لا فائدة منه إذا كانت الطائرة لا ترى المستخدم، والتتبع وحده لا يوفر وسيلة للتحكم. يهدف هذا المشروع إلى دمج هاتين الوظيفتين من خلال تطوير نظام ذكي يعمل فيه نموذجان معاً: نموذج للتعرف على الإيماءات وآخر لكشف الأشخاص. الهدف هو إنشاء واجهة تفاعل إنسانية-آلية سلسلة وقوية وغير مقيدة، حيث يصبح جسد المستخدم هو أداة التحكم الفعلية.

الكلمات الأساسية:

تفاعل الإنسان مع الطائرة بدون طيار، التعرف على الإيماءات، تتبع الأشخاص، الرؤية الحاسوبية، التحكم في الطائرات بدون طيار، التفاعل الطبيعي.

Dedication:

In the name of Allah, the Most Gracious, the Most Merciful, without whom nothing would have been possible, I dedicate this work,

To myself. For never giving up in the face of trials, for drawing from every difficulty the strength to carry on. This thesis is the culmination of an arduous path, and proof that perseverance always triumphs in the end.

To my dearest parents, my foundation and my source of inspiration. For your unconditional love, your silent prayers, your sacrifices, and your unwavering faith in my abilities. This thesis is the fruit of your steadfast support.

*To my beloved sisters, **Manel, Nesrine, Ferial, and Chahrazed**, and to my brother **Lyes**. Your presence, your encouragement, and our bond have been precious treasures throughout this journey.*

*A very special dedication goes to my little brother, **Zakaria**. Your joy for life and your affection were rays of sunshine in the darkest moments.*

*To my precious friends, **Wiam, Naila, Ghazlene, Rihab, Kheira and suka**, as well as to all my friends. For the moments of doubt turned into laughter, for your moral support, and for our sincere friendship which has been a true refuge.*

My deepest gratitude goes to my supervisors, for their patience, guidance, and invaluable advice. I also thank all the professors who contributed to my education and taught me so much.

Finally, from the bottom of my heart, I dedicate this success to my nephews and nieces. To me, you are like my own children, my greatest source of happiness and my daily motivation.

ACKNOWLEDGEMENTS

As I conclude this dissertation, I wish to express my deepest gratitude to all the individuals who, directly or indirectly, have contributed to its development and completion.

*My most sincere thanks go first and foremost to my supervisors, **Mr. Choutri** and **Mr. Lagha**. Their guidance, availability, and insightful advice have been pillars of support throughout this project. I am grateful for their scientific expertise, their rigor, as well as for the trust they placed in me. Their pertinent feedback and encouragement were essential in overcoming challenges and in bringing this research to a successful conclusion.*

*I would also like to thank the **members of the jury** for the honor of accepting to evaluate this dissertation, and for the time they will dedicate to its reading. Your feedback will be invaluable.*

*I extend my gratitude to the entire **faculty and teaching staff** of institute of aeronautics and spatial studies Blida-1 for the quality of the education I have received over these years.*

*I do not forget my **classmates and friends**, with whom I have shared these years of study. The discussions, mutual support, late nights of work, and moments of relaxation greatly contributed to creating a stimulating and pleasant environment. Thank you for your mutual support and your friendship.*

*Finally, my most heartfelt thanks go to **my family**. Your love, your infinite patience, and your unconditional moral and material support have been my greatest strength. Without your constant encouragement and your confidence in me, this project simply would not have come to fruition. May this work stand as a testament to my eternal gratitude.*

Table of content

Abstract:	12
Dedication:	14
ACKNOWLEDGEMENTS	15
Table of content	16
list of figures.....	19
list of tables	20
General introduction.....	22
Chapter 01: Techniques for Human-Drone Interaction.....	14
Introduction:	14
I.1.Definition of HDI.....	15
1.2. Research on Human-Drone Interactions	16
1.2.1. Role of humans in HD.....	16
1.2.2. HDI Research Over Time	17
I.2.3. Innovative Control Interfaces	17
I.2.4. User interfaces	18
I.3.4.5 Innovative Prototypes and Use Cases	30
2. Machine learning.....	35
2.1 Introduction to Machine Learning.....	35
2.2 How Machine Learning Works	36
2.2.1 Data Collection and Preparation	36
2.2.2 Algorithm Selection.....	36
2.2.3 Training the Model	37
2.2.4 Model Evaluation and Deployment.....	37
2.3 Types of Machine Learning.....	37
2.3.1 Supervised Learning.....	37
2.3.2 Unsupervised Learning	38
2.3.3 Semi-Supervised Learning	39
2.3.4 Reinforcement Learning.....	39
2.4 Common Machine Learning Algorithms.....	40

2.5 Machine Learning Applications	41
2.5.1 Everyday Use Cases	41
2.5.2 Advanced Applications.....	41
2.6 Benefits and Limitations of Machine Learning.....	41
2.6.1 Advantages	41
2.6.2 Challenges	41
2.7 Difference Between AI and Machine Learning.....	42
2.8 Machine Learning in Data Science	42
Conclusion.....	42
Chapter 02: Gesture Control	44
Introduction	44
2.1 Overview of HGR System Components	44
2.1.1 Gesture Identifiers	45
2.1.4 Contribution of the Paper	Error! Bookmark not defined.
2.1.4 Governing Criteria for System Design.....	46
2.2 Gesture Description Model	46
2.2.1 Overview	46
2.2.2 Gesture Type Justification	46
2.2.3 Gesture Model Selection Justification.....	47
2.2.4 Gesture Information Justification.....	47
2.3 Selection of Data-Acquisition Method.....	47
2.4 Selection of Gesture-Identification Algorithm.....	47
2.5 Validation of Selected Gesture-Identification Algorithm.....	47
2.6 Stage Five: Selection of Gesture Classification Algorithm.....	48
2.7 Stage Six: Gesture Mapping and Tuning.....	48
2.8. Results and Analysis.....	48
2.8.1 Training Process and Convergence	49
2.8.2 Quantitative Performance Analysis	50
2.8.3 Qualitative Analysis and Error Examination.....	52
2.9 Detailed Training Results	61
Chapter 03: Tracking of a Person.....	67

Introduction	67
3.1 Development Methodology	67
3.2 Results and Performance Analysis	69
3.3 Detailed Training Results	76
Conclusion.....	80
Chapter 4: System Architecture, Practical Implementation, and Testing.....	82
Introduction	82
4.1 Hardware Architecture of the Prototype.....	82
4.2 Prototype	85
4.3 Software Environment Configuration	86
4.3 Control Algorithm and Integration Logic.....	87
4.5 Testing, Calibration, and Refinement Protocol	89
Conclusion.....	90
General conclusion:	91

list of figures

Figure 1- 1: The four major fields of Human drone interaction research	15
Figure 1- 2: Sample screen of the graphical user interface developed for the Aero stack.	20
Figure 1- 3: The Dynamics Viewer used in the GUI.....	21
Figure 1- 4: Description of High-level setup for the proposed body position NUI.	23
Figure 1- 5: High-Leved Description of Visual Marker NUI.....	24
Figure 1- 6: From left to right, the stop, continue, and come hand gesture data labelled with hand landmarks	25
Figure 1- 7: High-Level Description of Speech Command NUI	26
Figure 1- 8: Flying drones with brains	27
Figure 1- 9: Medical drone and Figure 1- 10: A drone against forest fires	32
Figure 1- 11: Life-saving drones	32
Figure 1- 12: Drones for the blind.....	34
Figure 1- 13: machine learning techniques	36
Figure 1- 14: Types of machine learning.....	37
Figure 1- 15: Supervised learning	38
Figure 1- 16: Unsupervised learning.....	39
Figure 1- 17: Reinforcement learning	40
Figure 2- 1: Training and validation metrics over 50 epochs.....	49
Figure 2- 2: Precision-Recall curve for all gesture classes.	50
Figure 2- 3:F1-Confidence curve, showing an optimal F1 score of 0.99 at a confidence threshold of 0.430.....	52
Figure 2- 4: <i>Normalized Confusion Matrix</i>	53
Figure 2- 5: Example predictions from a validation batch.....	54
Figure 3- 1: Training and validation metrics over 50 epochs.....	69
Figure 3- 2: Precision-Recall curve for the "person" class.	71
Figure 3- 3: F1-Confidence curve, indicating an optimal F1 score at a confidence threshold of 0.755.....	71
Figure 3- 4: Normalized Confusion Matrix.....	72
Figure 3- 5: Example detections on augmented validation images.....	73
Figure 4- 1: Photograph of the Assembled.....	84
Figure 4- 2: Electronic Wiring Diagram (ref a suka men fritizing)	86
Figure 4- 3: Flowchart of the Main Control Loop	87
Figure 4- 4: Block Diagram of the Closed-Loop Control System	88

list of tables

Table 1- 1: Roles of Humans in HDI.....	16
Table 1- 2: Control Interface Summary.....	18
table 2- 1: Per-Class Performance Metrics on the Validation Set	51
table 2- 2: Results.....	61
table 2- 3: Results.....	62
Table 3- 1: Results.....	77
Table 3- 2: Results.....	78

Abbreviation:

UAVs - Unmanned Aerial Vehicles
HDI - Human-Drone Interaction
HRI - Human-Robot Interaction
NUI - Natural User Interface
BCI - Brain-Computer Interface
EEG - Electroencephalography
fNIRS - Functional Near-Infrared Spectroscopy
GUI - Graphical User Interface
HGR - Hand-Gesture Recognition
YOLO - You Only Look Once
CNN - Convolutional Neural Network
mAP - mean Average Precision
IoU - Intersection over Union
PWM - Pulse Width Modulation
CSI - Camera Serial Interface
GPIO - General Purpose Input/Output
GCS - Ground Control Station
BLOS - Beyond-Line-Of-Sight
VR - Virtual Reality
SVM - Support Vector Machine
KNN - K-Nearest Neighbors
ML - Machine Learning
AI - Artificial Intelligence
NLP - Natural Language Processing
FAA - Federal Aviation Administration

General introduction

Drones, sometimes referred to as unmanned aerial vehicles (UAVs), are robots that can fly on their own or with the help of various control methods like joysticks, smartphones, the human brain, voice, gestures, and more. Drones were complicated machines that were only accessible to the military until the early 2000s.

The development of smaller, more manageable, and less expensive systems is made possible by recent developments in hardware and software technology. Nowadays, drones are used for a wide variety of civilian tasks and their use is only anticipated to grow in the near future.

The FAA estimates that by 2022, there may be 3.8 million drones registered in its database. The adoption of natural user interfaces (NUIs), such voice commands and body movements, which have been tested in the state of the art, may further improve this. The majority of findings seem to suggest that using a NUI makes it easier for people to engage with drones.

This effort aims to answer the following question: How can a UAV be controlled and operated using gestures rather than joysticks? And is it possible to operate the drone using both gestures and speech?

For this, we were particularly interested in the design of a Multi model (speech and gesture) control of a quadrotors UAV

In the first chapter, we discuss human-drone interaction to help readers grasp the importance of drones in our daily lives and how to interact with them.

The second chapter discusses some fundamental ideas in voice recognition as well as the deep neural network, which has demonstrated notable advancements in the extraction and recognition of speech features.

The third chapter covers gesture command and the procedures for employing natural hand gesture movements to control the behavior of the UAVs.

The combination of input modalities, including speech and gesture data, is covered in the fourth chapter. the Presentation of the obtained experimental results is followed by an analysis of the performance (accuracy) attained.

In the forth chapter, we look at the implementation side, or how we applied the simulation to quadrotors while displaying the feedback in a graphical interface.

We wrap up our work with a summary of the results and recommendations for more research to preserve the performance and continuity of the suggested subject.

CHAPTER 01



Chapter 01: Techniques for Human-Drone Interaction.

Introduction:

In drone systems, humans play different roles depending on the drone's purpose and level of autonomy.

They may act as active controllers, directly piloting drones for tasks like racing or photography. Alternatively, they can be recipients, interacting with the drone without controlling it such as receiving deliveries or viewing drone advertisements.

Social interaction is also possible, where drones like *Joggobot* accompany users for engagement. With autonomous drones, humans often act as supervisors, either programming missions or monitoring operations to ensure safety.

Even in autonomous use, humans may still be involved indirectly as recipients.

In recent years, the field of Human-Robot Interaction (HRI) has expanded to include Human-Drone Interaction (HDI), driven by growing interest in how people interact with unmanned aerial vehicles (UAVs).

Originally developed for military purposes after World War I, drones are now widely used in civilian applications such as entertainment, delivery, agriculture, assistance for people with disabilities, sports, and rescue missions.

The rapid rise in drone usage brings both opportunities and challenges. While they enable innovative services, they also introduce potential risks to society. As a result, studying HDI is essential to better understand and improve human-drone collaboration across different domains.

I.1. Definition of HDI

Human-Drone Interaction (HDI) is an emerging field focused on understanding, designing, and evaluating drones for use by or with humans.

While it shares similarities with Human-Robot Interaction (HRI), HDI is distinct due to drones' unique ability to move freely in 3D space and their specific design characteristics.

The field is multidisciplinary: some researchers work on making drones more approachable and user-friendly, while others focus on developing intuitive control interfaces for inexperienced users. (1)

As drones become more common in everyday life, studying how humans interact with them is increasingly important.

This research aims to outline the current state-of-the-art in HDI. Key areas include:

- Developing new control methods
- Enhancing human-drone communication
- Measuring interaction distances
- Exploring new applications for drones



Figure 1- 1: The four major fields of Human drone interaction research

1.2. Research on Human-Drone Interactions

1.2.1. Role of humans in HD

Humans interact with drones in various ways, depending on the drone's function and its level of autonomy.

In some cases, the human acts as an active controller, directly piloting the drone through an interface to accomplish specific tasks, such as taking landscape photos or competing in drone races.

In other scenarios, the person may simply be a recipient, benefiting from the drone's presence without controlling it—like receiving a package or encountering a drone displaying advertisements in public. Interaction can also take on a social dimension, where drones serve as companions; for example, a drone like Joggobot flies alongside joggers, creating a form of social engagement.

With autonomous drones, the human role shifts to that of a supervisor, responsible for pre-programming flight missions or monitoring autonomous operations in case of unexpected events.

Even when drones operate independently, humans often remain involved in the loop, particularly as recipients of services or information.

These roles of humans in HDI are summarized in Table I.1 :

Table 1- 1: Roles of Humans in HDI

Role	Example
Active controller	Pilot controlling a drone during a photoshoot or racing competition
Recipient	User watching advertisement displayed by a drone or receiving a package delivery
Social Companion	User performing physical activities(e.g.jogging) in companion of a drone
Supervisor	Pilot supervising a wildlife monitoring autonomous mission

1.2.2. HDI Research Over Time

Human-Drone Interaction (HDI) is a relatively recent area of study within the research community. An analysis using Google Scholar, filtered by publication year, illustrates how interest in this field has grown over time. Before 2014, only two publications were found using the term “human-drone interaction.” However, by extending the search to include works published up to 2018, the number of results increases to 180. Although other search combinations involving terms like “drones,” “UAV,” and “human-in-the-loop” were explored, they mostly yielded results related to military applications and ethical issues. As a result, only the term “human-drone interaction” was used to track the evolution of research in this specific field.

1.2.3. Innovative Control Interfaces

Due to their complicated user interfaces, drones were first utilized for military purposes and required highly skilled pilots to operate. Researchers began to move interface design toward more contemporary user interfaces that no longer restrict drone operation to a remote controller or a ground control station as drone technology became more widely available and reasonably priced. Users may engage with drones through gesture, speech, sight, touch, and even brain-computer interfaces (BCIs) like electroencephalography (EEG) thanks to these cutting-edge techniques, sometimes referred to as natural user interfaces (NUI). In this part, we discuss novel control modalities and their effects on the pilot's experience. Each control interface affects the pilot's contact with the drone in a number of ways, including training duration, accuracy, latency, and interaction distance. A summary of the advantages and disadvantages of each NUI, as well as the traditional remote-controller interface, can be seen in Table I.2. Each NUI is further discussed in the sub-section below.(3)

Table 1- 2: Control Interface Summary

Control Interface	Advantages	Disadvantages
Remote Controller	Low-latency Precise control	Less intuitive than natural user interfaces
Gesture	Natural/intuitive Most drones already carry a camera Short training period	Collocated interaction only High latency and lower control precision when compared to others
Speech	Natural/intuitive Short training period Handless	Propeller sound can impact speech recognition Not all drones already have microphones Collocated interaction only
BCI	Handless Accessible to users with physical disabilities Ability to measure pilot's Affective and cognitive state	Longer training periods Lower control accuracy When compared to others Slower input and output response
Touch	Natural/intuitive short training period	Collocated interaction only requires safety measures
Multimodal	Can combine benefits from various modalities	Possible higher complexity and cost

I.2.4. User interfaces

In this section the implemented user interfaces will be explained. The point of interaction between a computer, website, or application and people is known as the user interface (UI).

A good user interface (UI) should make the user's experience as easy and uncomplicated as possible, requiring the least amount of effort on the part of the user to achieve the greatest number of desired outcomes. The GUI and NUI formats are the most crucial.(2)

First of all, we will explain GUIs, which have a big impact on HMI and are a great tool for robotics, thus the need to implement a robust and reliable interface for the drone platforms which will be explained. then we will talk about NUI which can provide efficient communication.

I.2.4.1. Graphical User Interface (GUI) :

This section describes how the operator can interact with quadrotor platforms using the developed GUI.

In general, a GUI provides some functions to help operators in certain tasks that are difficult to be supported by a NUI. For example, they correspond to tasks when the operator requires detailed information such as vehicle set up or mission monitoring at software level (e.g., during software maintenance).(4)

The GUI allows the interaction with the vehicle, observing the states and dynamics and presents graphical views and images to help the user to understand both the external and internal behavior of the vehicle.

In general, the operator can use a GUI to perform the following types of tasks:

- Specify drone behavior in advance (vehicle set up)
- Monitor drone behavior during a mission
- Operate manually with simple movements
- Collect data for later use

Fig. I.2 shows a sample screen of the user interface with the windows layout with main parts: the control panel on the upper left side of the image, the dynamics viewer on the lower left side, the windows for detailed content on the righthand side of the screen with different tabs (parameters, the camera viewer or the performance monitor). At the top, there are drop-down menus (file, view, settings, etc.) to perform additional tasks

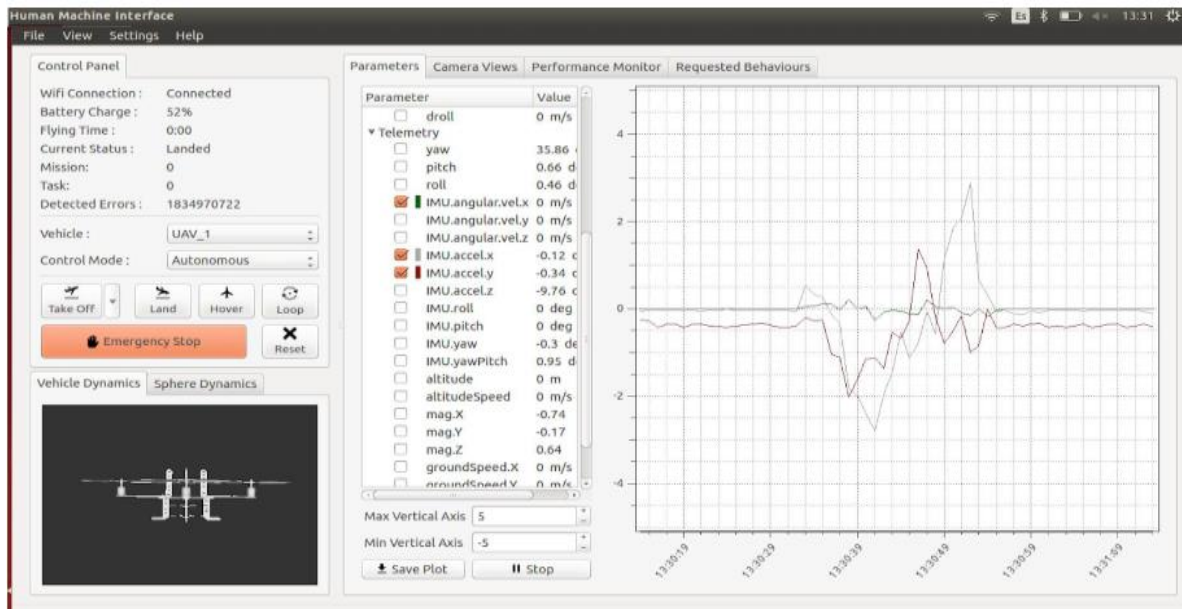
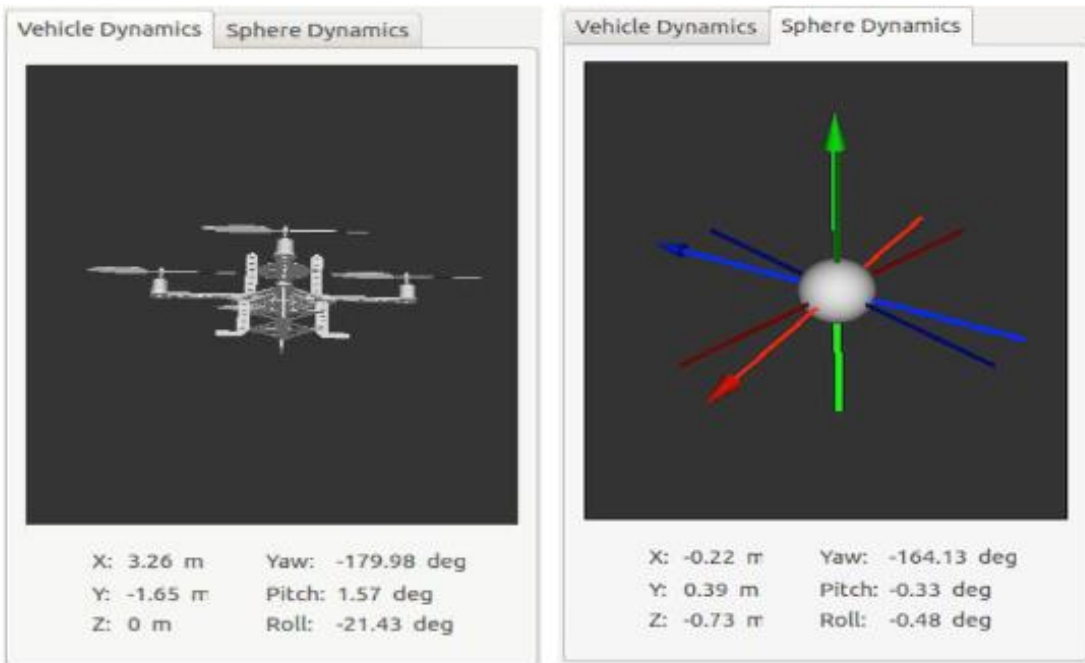


Figure 1- 2: Sample screen of the graphical user interface developed for the Aero stack.

The following sections describe the different parts of the graphical user interface.

- The control panel: It displays the primary control commands as well as the overall condition of the system. This panel offers a synopsis of the crucial information that must always be shown on the screen.
- Parameter viewer: The operator may inspect the values of numerical parameters and view various plots of the parameters chosen in real time to examine and contrast their values, allowing them to keep a close eye on the vehicle's behavior during a mission.
- Camera viewer: This tool displays photographs and/or videos that were taken by the aerial vehicle while it was in flight.
- Requested Behaviors Viewer: The operator can consult and request the activation of particular behaviors via the behavior viewer. The GUI displays a list of the vehicle's possible behaviors along with the condition of each one.



(a) The Vehicle Dynamics

(b) The sphere Dynamics.

Figure 1- 3: The Dynamics Viewer used in the GUI

- This viewer presents two animated images as seen in Fig.I.4: (1) the vehicle dynamics, a 3D representation of the vehicle (Fig. I.4(a)) and (2) the sphere dynamics, a sphere with orientation axis (Fig. I.4(b)). In the sphere representation, there are three fixed axis that represent the reference system, and three variable axis that represent the orientation changes.

I.2.4.2. Natural User Interfaces (NUIs)

However, compared to a traditional GUI, a NUI can offer more effective communication for some operator responsibilities. Voice instructions, which can be utilized more effectively in conjunction with other communication modalities, or gestures, which are simpler to learn, can be used to steer the car during manual operation. In general, human-drone cooperative work for some complicated tasks in dynamic circumstances where human judgments might supplement the

incomplete information used by the vehicle can benefit greatly from a communication based on NUI. (5)

In the following sections a description of the types of NUIs implemented and the setup used for each interaction will be given. First, the vision based NUIs will be explained; starting with body and followed by marker interaction, hand and speech interaction.

➔ Visual Body Interaction

The drone, equipped with an on-board camera, uses computer vision algorithms to detect a person and track it in the image plane.

With the previous knowledge of the approximate dimensions of the object tracked, the drone is able to reconstruct the 3D relative position of the object with respect to the drone.

A control algorithm sends commands to the drone ensuring it maintains the distance (x, y and z) and point of view (yaw angle).



Figure 1- 4: Description of High-level setup for the proposed body position NUI.

The high-level description of the suggested NUI is shown in Fig. I.4. This controller will be used to investigate a novel kind of human-drone interaction in the context of the work that is being presented.

The computer vision algorithm is used for this purpose in order to identify and follow the individual it is dealing with. Even if the person runs or moves, the drone can still track them, maintaining its distance and field of vision. In order to safeguard the user while in flight, the quadrotor is always used with an indoor hull. One of the first animal behaviors is displayed by this NUI.

➔ Visual Marker Interaction

Many animals, including humans, use color, depth, and motion as visual signals. Since visual cues and markers employ the camera, which is likely the most crucial sensor in robotics or drones, they are a widespread technique in robotics.

For dependable and precise target tracking or detection, basic monochromatic cameras might be employed. In this kind of human-drone interaction, the user instructs the drone on what to perform by manipulating markers, which are also known as tags in the literature.

The interface is simple and the user feels included in the decision-making process while controlling the drone because no other equipment is required save the on-board front-facing camera. Robust and accurate interaction is achieved while preserving a high degree of user ease through the use of carefully designed markers.



Figure 1- 5: High-Leved Description of Visual Marker NUI

→ **Hand Gesture Interaction**

In this work, a gesture recognition system is implemented to enable human-drone interaction using six predefined hand gestures: **up**, **down**, **left**, **right**, **on**, and **off**. The dataset was created by recording short gesture videos using a standard webcam. Frames were then extracted and manually annotated using **Label Studio**, an open-source data labeling tool.

For feature extraction, **MediaPipe** is employed to detect and track hand landmarks. MediaPipe's hand tracking pipeline consists of a palm detection model and a hand landmarks model, which together provide accurate 3D coordinates of 21 hand keypoints. These features are then used to classify the gestures.

The approach focuses on ensuring robustness and generalization by capturing gesture variations across different hand positions, lighting conditions, and orientations. This process enables the system to recognize commands reliably and serves as a foundation for hands-free control of the UAV in a future deployment.

In the feature extraction phase, MediaPipe is used to extract hand landmarks. As illustrated in Fig. I.6, MediaPipe applies a robust model based on CNN to determine the keypoint localization of 21 hand-knuckle coordinates inside the detected hand regions. Mediapipe model was trained using around 30K realworld images as well as synthetically created hand models with a variety of backgrounds. A palm detection model and a hand landmarks detection model are included in the MediaPipe hand landmarker model bundle. The palm detection model detects hands inside the input image, while the hand landmarks recognition model recognizes specific hand landmarks on the palm detection model's cropped hand image.



Figure 1- 6: From left to right, the stop, continue, and come hand gesture data labelled with hand landmarks

Several studies on gesture-based control systems highlight their intuitive and user-friendly nature, making them accessible with minimal training requirements. One of the main advantages of this interaction modality is its device-free operation, eliminating the need for physical controllers such as joysticks. However, despite these benefits, gesture control may not be the most suitable solution for applications requiring high precision and fine control, due to its relatively higher latency and lower accuracy compared to more traditional input methods

➔ **speech command Interaction**

The increasing number of easily accessible voice recognition software programs and toolkits has led to an increase in the frequency of this kind of contact for HCIs. voice-to-text applications were the primary use case for early voice recognition interfaces. Products like Dragon Naturally Speaking enable the user to send instructions that the software recognizes as such, have a document synthesized as an audio stream, or dictate and have voice transcribed as written text. These programs provide hands-free, user-friendly interfaces that bridge the gap between spoken and written language. Because it allows the user to focus all of their visual attention on giving commands to the drone, this kind of interaction can improve the user experience in drone applications.

Fig. I.7 shows the high-level setup of the proposed speech command NUI. As can be seen, the setup is simple. No devices other than a microphone are needed to command tasks to the drone.

Voice commands are sent to a Ground Station to be processed and later sent as tasks to the drone using the Aerostack.



Figure 1- 7: High-Level Description of Speech Command NUI

→ Brain-Computer Interaction (BCI)

The ability of Brain-Computer Interface (BCI) technology, especially non-invasive technologies, to revolutionize human-computer interaction (HCI) has made it a focal point of innovation. Acquiring neural signals from the brain, decoding the signals to derive basic goals, and evaluating these intentions to convert them into various computer output formats are the three primary stages it takes to connect human brain signals to direct computer commands.(6)

Non-invasive BCIs use external sensors, such as EEG or fNIRS, to record brain activity painlessly and without the need for surgery, in contrast to invasive BCIs that need surgical implants. By giving persons with physical limitations an other way to engage with the outside world and by giving the general public additional methods to communicate with technology, this method increases the accessibility of BCIs.

By enhancing user experience and bringing in new paradigms of computer interaction, the development of non-invasive brain-computer interfaces (BCIs) shows great promise to open up new possibilities in human-computer interaction (HCI). The technical developments in non-

invasive BCI systems that influence HCI are examined in this article, with particular attention paid to important elements including signal collection, brain decoding models, and machine learning algorithms that help to maximize BCI performance.

Examining how developments in non-invasive BCI systems improve HCI and investigating the theoretical underpinnings of the models and methods that support the precision and effectiveness of BCI-driven interactions are the goals of this study.



Figure 1- 8: Flying drones with brains

→ **Multimodal**

Integrating different interaction methods combines the strengths of each. Studies show that many users naturally adopt multimodal interaction with drones. For instance, drones can be controlled by voice commands for takeoff and by hand gestures for directional movement. Some prototypes use onboard sensors to detect both speech and gestures, eliminating the need for external devices. Voice commands are used for takeoff, while a whistling sound—less affected by propeller noise—is used for landing. Movement is guided by tracking hand gestures using colored gloves.

Virtual reality environments have also been used to explore multimodal control of drone swarms in simulated search-and-rescue missions. Users controlled multiple drones using speech and gestures. Experiments highlighted challenges like depth perception and emphasized the need for clear commands, immediate system feedback, and enriched visual interfaces such as artificial shadows under drones to improve spatial awareness.

Three main information flows exist between users and ground control stations: (1) commands from the operator to the ground control station (GCS), (2) system feedback from the GCS to the user, and (3) physiological data from the user (e.g., heart rate or EEG) used to adapt the system interface. Multimodal interaction enhances these exchanges by increasing communication channels, managing information overload, and adapting to diverse environments.

➔ Other Control Interfaces

Recent developments in human-drone interaction have focused on making control more natural, immersive, and intuitive. Some commercial drones, like those from DJI and Ryze, are designed to be “safe-to-touch,” enabling physical interactions such as hand landings. Studies show that 58% of users prefer this mode of control when propeller guards are present, citing reduced mental workload compared to traditional joysticks. Gaze tracking has also been explored, allowing users to control drone yaw and pitch with eye movements, while other controls are handled by the keyboard, resulting in effective 3D navigation.

The **Birdly** system offers an even deeper level of immersion, where users lie prone and mimic bird-like flight by controlling a real drone with body movements—arms for roll and pitch, and head for camera orientation—enhanced by motion platforms, VR goggles, and wind feedback. Users found it more enjoyable and immersive than joystick control. Another approach,

Flying Head, links head movements directly to the drone's position and orientation, enabling intuitive control that outperformed traditional joysticks in terms of speed, accuracy, and task efficiency during target-tracking scenarios. This could lead to new forms of aerial sports based on full-body control. Finally, a physical 3D-printed map has been used as an interface for drone path planning: users draw flight paths directly on the model using a stylus, and the drone autonomously follows the trajectory. Augmented reality is used to visualize the drone's path, showing the promise of tangible, spatial interfaces in drone control systems.

I.3.4.3. Distance, Communication and Emotion Encoding

→ Interaction distance

The interaction distance between the drone and the human must be taken into account for a positive social relationship. [8] In the prior experience, 47% of US users remained in the personal area (1.2m), 37% remained in the drone's intimate space (45cm), and the remaining sixteen percent engaged in social space (3.7 m). Nonetheless, 50% of the Chinese participants interacted in the intimate zone, 38% in the personal space, and 6% in the social space, indicating more relaxed and intimate interactions. Drones approached people at two different heights (1.80m and 2.13m) in another survey of users, and the results showed that height had no discernible effect on the acceptable approach distance.

→ Drone feedback

In order to prevent system misunderstandings that might possibly result in accidents, studies have previously looked into ways to identify the mutual attention between a drone and its users as well as successful communication.

A comparison of four distinct drone identification gestures was the focus of [10].

Users prefer rotation in the yaw axis to signify recognition, according to the data.

In [11], the capacity of a drone to communicate its purpose to people was examined. The expression involved the input and output of speed profiles as well as the manipulation of primitive motions utilizing arc trajectories. The following tasks are tested with the participants by the drone prototypes constructed with these manipulations: approaching an object (anticipation), avoiding an

object (arc + easy entrance and exit), moving away from an object (arc + easy input and output), and getting close to a person (easy input and output). According to the findings, users would rather operate a drone with altered flight paths than with simple ones for safety reasons and for a more organic and intuitive interface.

➔ Remote Communication

Humans can interact with drones either locally or remotely. With advancements like 5G, remote control is now more reliable, enabling applications such as package delivery and drone swarms. Drones can even act as mobile network stations. However, beyond-line-of-sight (BLOS) flights raise safety and security concerns, including the risk of crashes and cyber-physical attacks like jamming or delivery delays.

➔ Emotion Encoding

Drones can express emotions through movement parameters such as speed, trajectory, and reaction time, despite lacking humanoid features. A study tested three emotional flight profiles (exhausted, antisocial, adventurer) by adjusting flight behavior. Participants observed the drone and identified its emotional state. Results showed emotions were correctly recognized 60% of the time with one keyword and up to 85% with an additional descriptive word. This highlights that emotion encoding can improve human-drone communication.

I.3.4.5 Innovative Prototypes and Use Cases

Although there are many uses for drones today, researchers are always looking for new methods to make these devices useful.(7)

➔ Flying User Interfaces

This section explores the use of drones as mobile and interactive user interface platforms. Thanks to their ability to move freely in 3D space, drones can position themselves dynamically around users and act as both input and output devices. Researchers have categorized flying

interface interactions into three phases: *approaching*, *interacting*, and *leaving*, allowing drones to engage users effectively in tasks such as public information dissemination, sports guidance, or even crowd control during emergencies.

Several prototypes have demonstrated innovative implementations. One project used two Parrot AR.Drones, one carrying a projector and the other a screen, to create a flying public display. Visual markers and computer vision allowed the projector drone to track the screen drone, ensuring accurate image alignment. Experiments showed that with modified stabilization algorithms and visual markers, the drone system achieved significantly better hovering stability and projection accuracy.

Another study employed an octocopter carrying a smartphone and projector to display SMS messages on building walls. During a live outdoor test, users reported that the system was fun, attention-grabbing, and suitable for applications like advertising and storytelling. A similar indoor prototype allowed event attendees to text messages, which were then displayed on a hovering drone. The audience reacted with enthusiasm, demonstrating the potential for drones to foster group engagement and social interaction.

Despite their effectiveness, some limitations were identified. A user study evaluated readability under different motion conditions. Results revealed that while users could comfortably read content from stationary or slowly moving drones, readability decreased significantly when both the user and the drone were moving. Consequently, content design (such as font size) must be adapted to account for motion dynamics.

In summary, drones present a novel and flexible approach to augmenting user interfaces, particularly for interactive and attention-based applications. Their mobility, visibility, and ability to adapt to different contexts make them promising tools for both indoor and outdoor human-computer interaction systems.



Figure 1- 9: Medical drone



Figure 1- 10: A drone against forest fires



Figure 1- 11: Life-saving drones

➔ Social companions

One of the early prototypes in this area is **Joggobot**, a drone-based jogging companion designed not to improve physical performance, but to enhance the **social and motivational aspects** of jogging. Built using a Parrot ARDrone, it follows the runner at a fixed distance by detecting a visual marker on the user's shirt. Though the system has limitations (e.g., straight-line path only), preliminary user feedback was encouraging. Some appreciated the distraction from fatigue and motivation boost, while others wished for more control over pace. Suggestions for future versions included integrating **heart rate monitoring** for adaptive pacing based on health data.

In parallel, technical advancements were made with **custom-built quadcopters** for jogging. One such prototype used high-performance components including Pixhawk flight controllers and GPS modules to ensure stability, safety, and efficient outdoor tracking.

Drones are also being explored as **assistive companions for visually impaired users**. A conceptual system envisions a drone docked on a wearable bracelet, ready to be deployed via **voice commands**. It would navigate to a user-defined location while the person follows using **auditory feedback** from the propellers. Although the system remains under development, a **Wizard of Oz** study with a blind participant confirmed its viability, as the user was able to follow a miniature drone and responded positively to the idea.

In the domain of **environmental awareness**, drones have been proposed as persuasive agents to promote **cleanliness and eco-friendly behavior**. One study simulated a drone encouraging users to pick up trash and guide them to bins. It tested different interaction modes (visual, audio, and combined) through video prototypes. Although the mode of persuasion didn't significantly affect compliance, the study found that **gender and cultural background** influenced perception: females and participants from developing countries were more receptive to the drone's behavior, indicating that **human-drone interaction is shaped by cultural and demographic factors**.



Figure 1- 12: Drones for the blind.

➔ Arts and Sports

This section highlights how drones can be creatively used for **art creation** and **sports enhancement**.

- **Landscape Art:** A user can draw a sketch on a smartphone showing the drone's live camera feed. While the drone hovers, the user marks the ground by walking along the sketch lines. For example, a large smiley face was created on a grass field in under 30 minutes by mowing along the marked path.
- **HoverBall:** A drone encased in a soft spherical cage is used as a ball whose physics can be modified—changing speed, gravity, or flight trajectory. This enables the creation of new sports or adapting existing games for people with different skill levels, such as children or the elderly.

These examples show that drones can go beyond technical applications and become **interactive tools for creative, playful, and inclusive experiences**.

➔ **Haptic Feedback for Virtual Reality**

This section highlights two projects using drones to deliver haptic feedback in virtual reality (VR). Since traditional VR lacks tactile sensations, small quadcopters were employed to simulate physical interactions. In one project, drones like the Parrot Rolling Spider, equipped with safety cages and object tips, flew toward users at various speeds to simulate environmental effects such as bees, arrows, or falling debris within a VR jungle setting. In another project, a drone carrying a flat object like paper allowed users to feel resistance through airflow, even without direct force. A user study with four participants confirmed that the drone-generated haptic feedback improved accuracy in virtual tasks like drawing lines. These approaches demonstrate the potential of drones to enhance immersive VR experiences by adding a physical touch component.

2. Machine learning

2.1 Introduction to Machine Learning

Machine Learning (ML) is a fundamental subfield of Artificial Intelligence (AI) focused on enabling machines to learn from data and improve their performance over time without being explicitly programmed. ML algorithms analyze and identify patterns in various types of data, such as text, numbers, images, and sensor readings, in order to perform tasks or make predictions.

The core idea of ML is to build models that can generalize from observed data to unseen situations. These models are trained using historical data (training datasets), and once trained, they can be applied to new data to make decisions or predictions.

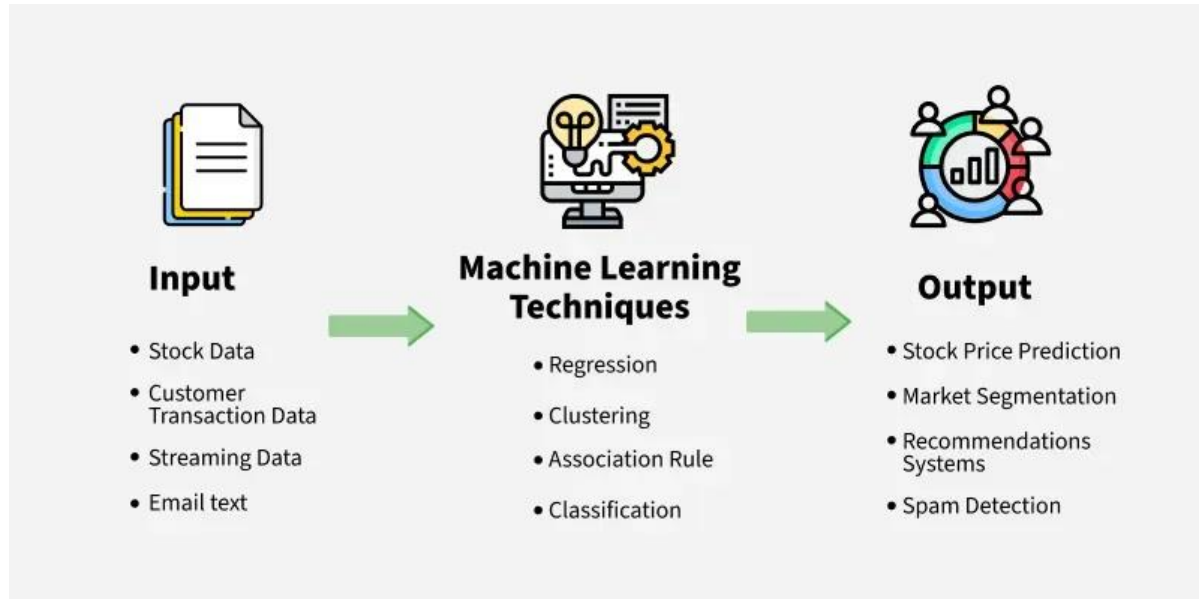


Figure 1- 13: machine learning techniques

2.2 How Machine Learning Works

The development of a machine learning model generally follows four key steps:

2.2.1 Data Collection and Preparation

The process begins with gathering and organizing a training dataset. The data can be labeled (supervised learning) or unlabeled (unsupervised learning), depending on the learning method. Data cleaning, formatting, and normalization are crucial to ensure accuracy and prevent bias during training.

2.2.2 Algorithm Selection

Choosing the right algorithm depends on the nature of the task (e.g., classification, regression, clustering) and the volume and structure of the dataset. Common algorithms include decision trees, support vector machines (SVM), logistic regression, and neural networks.

2.2.3 Training the Model

The algorithm is trained by feeding it data and adjusting internal parameters (weights and biases) to reduce prediction error. This process is iterative and involves comparing the model's output with the expected results.

2.2.4 Model Evaluation and Deployment

After training, the model is validated on new or test data. If it performs well, it can be deployed in real-world applications and continue to learn and adapt to new data through continuous feedback.

2.3 Types of Machine Learning

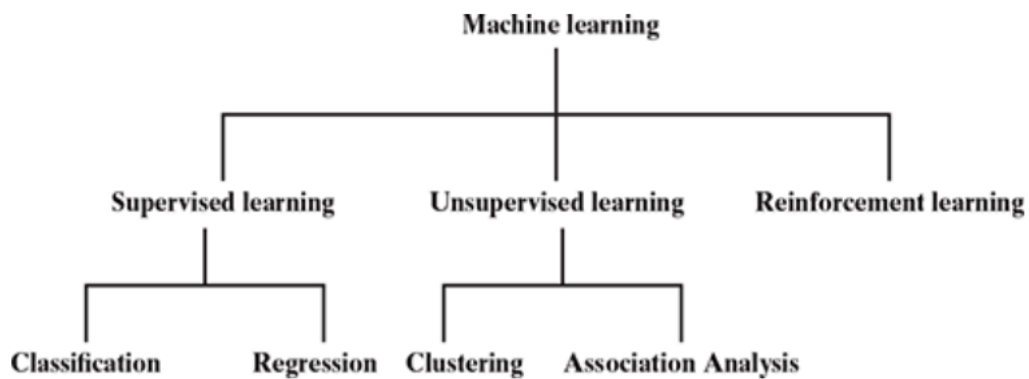


Figure 1- 14: Types of machine learning

2.3.1 Supervised Learning

In this method, the model is trained on labeled data, where both inputs and expected outputs are known. It is used for tasks like spam detection, email categorization, and medical diagnosis.

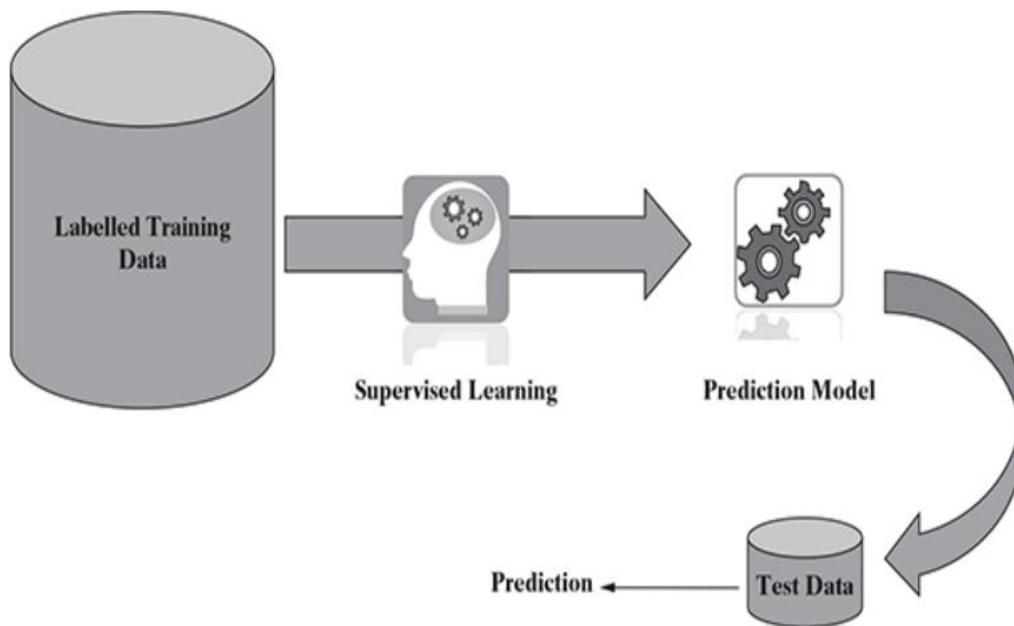


Figure 1- 15: Supervised learning

2.3.2 Unsupervised Learning

Here, the model explores the data without labeled outputs to identify hidden patterns or groupings.

It is commonly used in customer segmentation, anomaly detection, and market basket analysis.

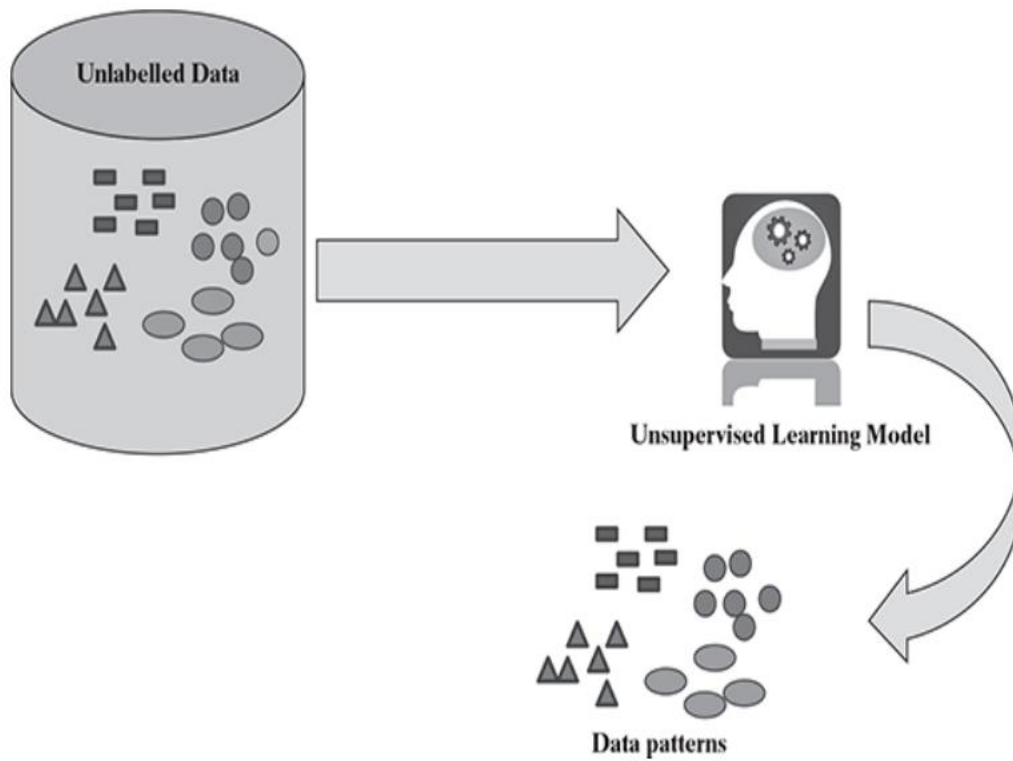


Figure 1- 16: Unsupervised learning

2.3.3 Semi-Supervised Learning

This technique combines a small amount of labeled data with a large volume of unlabeled data. It is useful when labeled data is expensive or time-consuming to obtain.

2.3.4 Reinforcement Learning

In reinforcement learning, the algorithm learns by interacting with an environment and receiving feedback through rewards or penalties. It is commonly used in robotics, gaming, and autonomous vehicles.

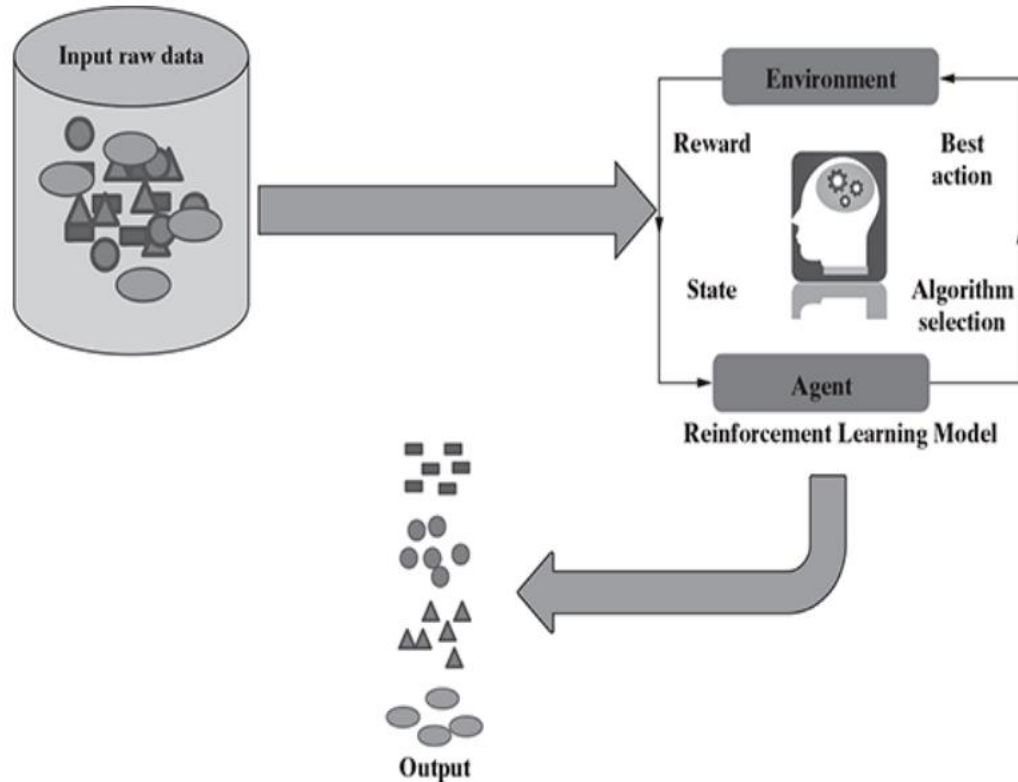


Figure 1- 17: Reinforcement learning

2.4 Common Machine Learning Algorithms

- **Linear Regression:** Predicts numerical outcomes based on the linear relationship between variables.
- **Logistic Regression:** Used for binary classification tasks.
- **Decision Trees:** A flowchart-like model that makes decisions based on data features.
- **Random Forests:** An ensemble of decision trees that improves prediction accuracy.
- **Support Vector Machines (SVM):** Effective in high-dimensional spaces for classification tasks.
- **Clustering (e.g., K-means):** Groups similar data points without prior labels.

- **Neural Networks:** Composed of layers of interconnected nodes that simulate the human brain, useful for image recognition, NLP, and complex pattern detection.

2.5 Machine Learning Applications

2.5.1 Everyday Use Cases

- **Speech Recognition:** Converts spoken language into text using NLP.
- **Customer Service:** Chatbots that handle FAQs and offer product recommendations.
- **Computer Vision:** Enables machines to interpret and act on visual data (e.g., facial recognition, self-driving cars).
- **Recommendation Systems:** Suggest products or content based on user behavior (e.g., Netflix, Spotify).

2.5.2 Advanced Applications

- **Medical Diagnosis:** Analyzing medical images to detect diseases.
- **Autonomous Vehicles:** Using sensors and learning algorithms to navigate.
- **Financial Forecasting:** Predicting market trends and risks.

2.6 Benefits and Limitations of Machine Learning

2.6.1 Advantages

- Automation of complex tasks
- Scalability across various domains
- Ability to discover hidden patterns in large datasets
- Continuous improvement with more data

2.6.2 Challenges

- Requires large, high-quality datasets
- Risk of algorithmic bias

- Interpretability of complex models like deep neural networks
- High computational costs

2.7 Difference Between AI and Machine Learning

Artificial Intelligence is a broad field focused on creating systems that mimic human intelligence. Machine Learning is a subset of AI that enables machines to learn from data. Other branches of AI include natural language processing, robotics, and expert systems. ML provides the statistical foundation for many modern AI applications.

2.8 Machine Learning in Data Science

Machine learning plays a central role in data science by automating data analysis, anomaly detection, clustering, and forecasting. It helps data scientists build predictive models that adapt over time and improve decision-making without direct human intervention.

Conclusion

Machine Learning has become a transformative technology with vast applications across industries. By enabling machines to learn from data, ML enhances automation, improves efficiency, and opens new possibilities in both research and practical fields. However, its success depends on the quality of data, ethical implementation, and ongoing human oversight.(10)

Chapter 02



Chapter 02: Gesture Control

Chapter 02: Gesture Control

Introduction

Alternative-control algorithms consist of two main components: a non-standard human–computer interface (HCI) and a command mapping algorithm [1–4]. These systems are deemed effective when the alternative HCI surpasses conventional control methods in terms of intuitiveness, flexibility, or accessibility. The effectiveness of such systems is generally evaluated using criteria such as recognition accuracy, ease of use without holding any device, short learning time, low implementation cost, and low computational complexity. The present project aims to develop such a system, where a **Raspberry Pi** manages a **Raspberry Pi camera** mounted on a **pan-tilt mechanism driven by two servo motors** to track and interpret a user's gestures in real-time.(14)

One prominent example of alternative HCI is hand-gesture recognition (HGR), where specific physical movements are detected and interpreted as control commands. Over the past decades, HGR has been extensively studied, evolving toward machine learning-based pipelines to improve precision and flexibility.

In recent developments, the **YOLO (You Only Look Once)** family of object detection models has gained wide popularity for real-time gesture identification thanks to its high detection accuracy, fast inference time, and lightweight architecture. YOLO was adopted as the core technology for this project, enabling the precise detection of hand gestures from RGB images with minimal delay. By training a YOLO model on a custom dataset that we created, we can achieve robust gesture recognition. This dataset was specifically designed for robustness, by capturing the gestures "**ON**," "**OFF**," "**UP**," "**DOWN**," "**LEFT**," and "**RIGHT**" with both the **left and right hands** at varying distances (**50 cm, 1 m, 2 m, 3 m, and 4 m**). The images were annotated using **Label Studio** and exported in the YOLO format for training.

2.1 Overview of HGR System Components

Modern hand-gesture recognition (HGR) systems can be decomposed into several key components: the data acquisition medium, the gesture descriptor, the gesture identification algorithm, and the gesture classification model [7–9].

Chapter 02: Gesture Control

Gesture descriptors are typically defined by three characteristics: (1) the physiological scope of the gesture (static vs. dynamic), (2) the type of information extracted (symbolic), and (3) the model representation used (bounding box). Symbolic gestures, which are most commonly used in real-time recognition systems, involve identifying hand shapes.

In our work, we use a **single RGB camera setup** (a PC webcam for data collection, a Raspberry Pi camera for deployment), combined with a **YOLO-based object detection model**, to detect symbolic hand gestures. This approach favors simplicity, cost-effectiveness, and avoids the need for specialized or wearable equipment. Moreover, the model is trained on annotated static gestures, which aligns with the symbolic gesture class and avoids the complexity associated with dynamic modeling.

2.1.1 Gesture Identifiers

Gesture identification refers to the method by which a human hand is detected apart from its background and transformed into a computer model used for classification.(11)

For this project, the **YOLO (You Only Look Once)** algorithm was selected to serve as both the gesture identifier and classifier. Unlike multi-stage pipelines that first detect a hand and then extract features (like MPH or InterHand2.6M which generate a 3D skeleton), YOLO operates as a single, end-to-end convolutional neural network (CNN). It processes an entire image in a single pass to directly output a set of **bounding boxes**, each associated with a confidence score and a class prediction. In our application, each bounding box represents an identified hand gesture, and the class prediction (e.g., "UP," "DOWN") serves as the recognized command. This unified approach is highly computationally efficient, making it ideal for real-time applications on resource-constrained hardware like the Raspberry Pi.(13)

2.1.3 Gesture Classifiers

Gesture classification refers to the process by which a feature extracted by the gesture-identification algorithm is classified as a particular gesture from a pre-defined list.

In this project's architecture, a **separate classification algorithm is not necessary**. The YOLO model is an object detector that integrates classification as a fundamental part of its detection pipeline. When it identifies a hand, it simultaneously assigns it a class label from the set

Chapter 02: Gesture Control

it was trained on (e.g., "LEFT," "RIGHT," "ON," "OFF"). This integrated approach greatly simplifies the overall system design, reduces computational overhead, and eliminates the need to select, train, and integrate a secondary classification model (such as an SVM, KNN, or Decision Tree). The classification is therefore handled directly by the final layers of the YOLO neural network.

2.1.4 Governing Criteria for System Design

To ensure the development of an effective and cohesive gesture-based alternative control system, a set of governing criteria was established. The selected criteria are listed below in order of importance:

1. Reliability in issuing the intended command
2. Reproducibility of gestures
3. Physically non-restrictive interaction
4. Ease of use and short learning curve
5. Low computational cost
6. Low monetary cost

2.2 Gesture Description Model

2.2.1 Overview

The first stage involved selecting a gesture description model. The final choices were:

- Gesture Type: Static, single-hand gestures
- Gesture Information: Symbolic
- Gesture Representation Model: 2D Bounding Box Model(15)

2.2.2 Gesture Type Justification

Static single-hand gestures were chosen to reduce algorithmic complexity. The gesture set (**ON, OFF, UP, DOWN, LEFT, RIGHT**) was designed to be visually distinct and intuitively map to drone commands, ensuring high recognizability and consistency.

Chapter 02: Gesture Control

2.2.3 Gesture Model Selection Justification

The **2D Bounding Box model** was selected as it is the native output of the YOLO object detection algorithm. This model was chosen to balance performance and computational efficiency. While more complex models like 3D skeletons provide more detailed information, they add significant computational overhead. For classifying distinct, symbolic static gestures, the shape and context captured within a 2D bounding box are sufficient for the YOLO model to achieve high accuracy. This choice avoids the need for complex landmark extraction and aligns perfectly with the goal of creating a fast, real-time system suitable for the Raspberry Pi.

2.2.4 Gesture Information Justification

Symbolic information was prioritized for recognizing discrete commands. This means each gesture corresponds directly to a predefined command, such as "Up" or "Land". While spatial information (like the position of the bounding box in the frame) is available and can be used for tracking by the servo motors, the primary classification relies on the symbolic meaning learned by the model from the image data.

2.3 Selection of Data-Acquisition Method

The goal was to identify a non-restrictive and cost-effective data-acquisition method. Based on the governing criteria, a **single RGB camera** was selected. Specifically, a **PC webcam** was used for data collection due to its convenience, and a **Raspberry Pi camera** is planned for the final deployment. This approach aligns with the criteria of low cost and an unencumbered interface.(12)

2.4 Selection of Gesture-Identification Algorithm

The **YOLO (You Only Look Once)** algorithm was directly selected as the core engine of the system. This choice was motivated by its state-of-the-art performance in real-time object detection, its unified architecture that simplifies the processing pipeline, and its efficiency, which makes it suitable for running on resource-constrained platforms like the Raspberry Pi.(17)

2.5 Validation of Selected Gesture-Identification Algorithm

This stage quantitatively validated the selected gesture-identification algorithm. The validation did not rely on clinical measurements but on **standard machine learning metrics**, computed on the validation dataset which the model had not seen during training. Performance was

Chapter 02: Gesture Control

evaluated using key indicators such as **Precision, Recall, and mean Average Precision (mAP)**. This process objectively verified the model's ability to correctly identify and classify gestures while also assessing its robustness to the variations in distance and hand introduced in the dataset.(18)

2.6 Stage Five: Selection of Gesture Classification Algorithm

This stage, as defined in a traditional pipeline, is **not applicable** to this project. As detailed in previous sections, the YOLO model performs gesture identification and classification simultaneously within a single neural network. There is therefore no need to select a separate classification algorithm.(19)

2.7 Stage Six: Gesture Mapping and Tuning

This stage translates recognized gestures into drone commands using a one-to-one mapping strategy.

- **Mapping Strategy:** Each of the six accurately classified gestures (**ON, OFF, UP, DOWN, LEFT, RIGHT**) will be assigned to a unique drone command.
- **Drone Platform:** The DJI Tello quadrotor is an ideal candidate platform for demonstration.
- **Tuning Parameters:** Three key parameters will be tuned experimentally: gesture hold time, command magnitude, and refresh rate, to achieve the smoothest and most intuitive control experience.

Excellent! This is exactly what is needed to write a complete and professional results section. Here is a structured proposal for presenting these results in your report. I have analyzed your graphs and extracted the key figures so that you can integrate them directly.(20)

2.8. Results and Analysis

This section presents the results obtained from the training and validation of the YOLO model. The analysis focuses on the model's quantitative performance, its training behavior, and a qualitative examination of its predictions.(6)

Chapter 02: Gesture Control

2.8.1 Training Process and Convergence

The model was trained for **50 epochs**. The training and validation progress was monitored to ensure proper learning and to prevent significant overfitting.

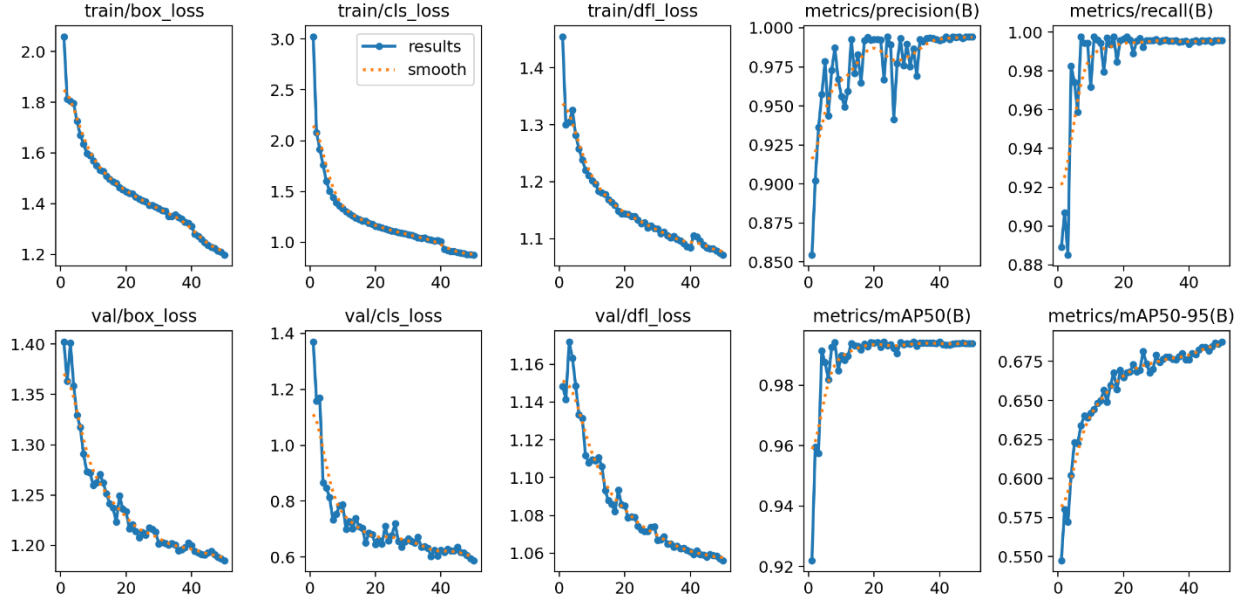


Figure 2- 1: Training and validation metrics over 50 epochs.

As shown in Figure 1, the loss functions (box_loss, cls_loss, dfl_loss) for both the training and validation sets show a consistent downward trend, stabilizing towards the final epochs. This convergence indicates that the model successfully learned to identify the features of the gestures from the dataset. Crucially, the validation loss (val/box_loss, val/cls_loss) does not show a significant upward divergence from the training loss, which suggests that the model generalizes well to unseen data without severe overfitting.

Concurrently, the key performance metrics on the validation set, such as Precision, Recall, and mAP, show a rapid increase in the initial epochs, followed by a plateau at a high level of performance. The metrics/mAP50(B) curve, in particular, converges to a value near **0.994**, confirming that the model reached a state of high accuracy and stability.

Chapter 02: Gesture Control

2.8.2 Quantitative Performance Analysis

The model's final performance was evaluated on the held-out validation set. The key performance indicators demonstrate an exceptionally high ability to accurately localize and classify the predefined gestures.

The primary performance metric, **mean Average Precision (mAP) at an IoU threshold of 0.5**, reached an outstanding **99.4%**. This confirms the model's excellent reliability. The Precision-Recall curve (Figure 2) illustrates this high performance across all classes.

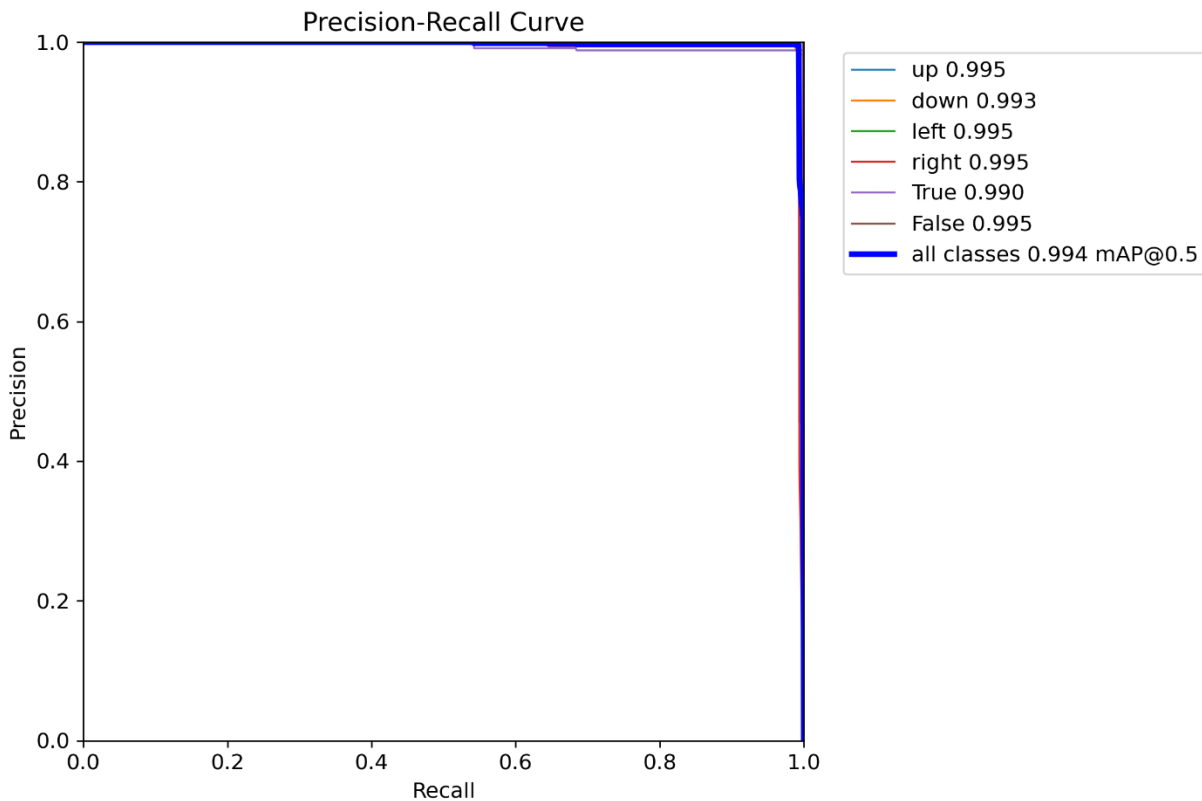


Figure 2- 2: Precision-Recall curve for all gesture classes.

Table 1 provides a detailed breakdown of the performance for each individual gesture class, based on the final mAP@0.5 scores shown in Figure 2.

Chapter 02: Gesture Control

table 2- 1: Per-Class Performance Metrics on the Validation Set

class	mAP@0.5
Up	0.995
Down	0.993
left	0.095
Right	0.995
True(on)	0.990
False (off)	0.995
All	0.994

Analysis: The results in Table 1 show that every single gesture class achieved a mAP score of 0.990 or higher. This indicates an extremely high level of reliability and satisfies **Criterion 1 (Reliability in Command Issuance)**. The model is equally proficient at recognizing all defined gestures.

The F1-Confidence curve (Figure 3) shows that the best balance between precision and recall is achieved at a confidence threshold of **0.430**, where the F1-score for all classes peaks at **0.99**. This is a critical insight for deploying the model, as it provides an optimal threshold for filtering detections.

Chapter 02: Gesture Control

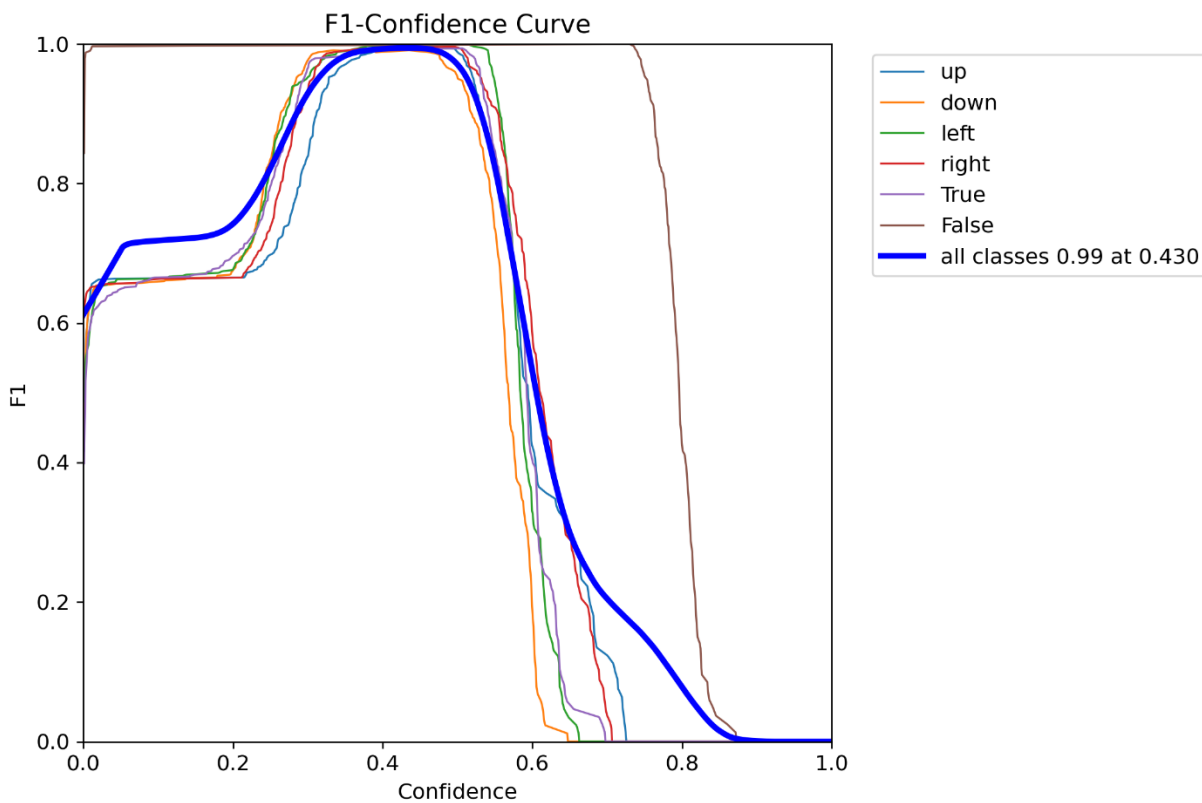


Figure 2- 3:F1-Confidence curve, showing an optimal F1 score of 0.99 at a confidence threshold of 0.430.

2.8.3 Qualitative Analysis and Error Examination

A qualitative analysis was conducted to understand the model's practical behavior.

Confusion Matrix:

Figure 4 presents the normalized confusion matrix. The strong diagonal values (close to 1.00) visually confirm the model's exceptional classification accuracy.

Chapter 02: Gesture Control

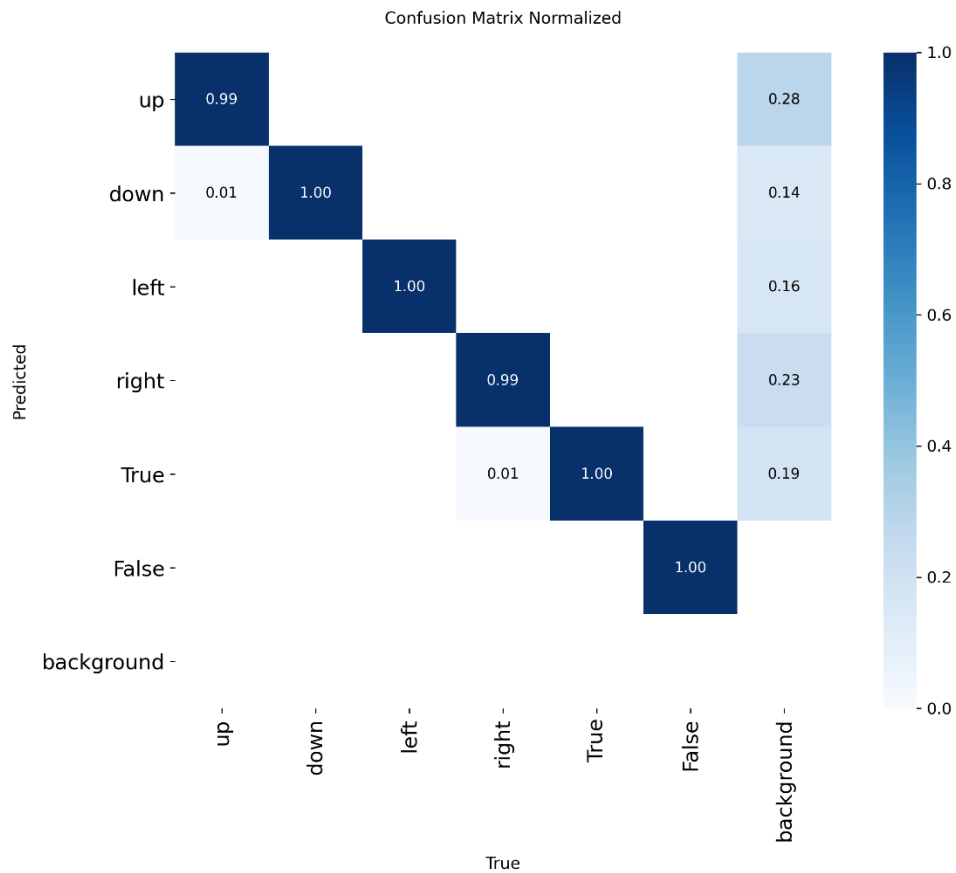


Figure 2- 4: *Normalized Confusion Matrix.*

Analysis: The confusion matrix reveals that misclassifications between defined gestures are nearly non-existent (e.g., only 1% of 'down' gestures were misclassified as 'up'). The main source of error is related to the background class.

- **False Negatives:** The background column indicates instances where the model failed to detect a gesture that was present. This occurred most frequently for the 'right' (23% of misses) and 'down' (14% of misses) gestures.
- **False Positives:** The background row (not shown in the normalized matrix but visible in the absolute one) indicates instances where the model incorrectly detected a gesture in the background. This was the single largest source of error, with 129 false positive detections.

Chapter 02: Gesture Control

Visual Examples of Performance:

The model's performance was further analyzed by examining its predictions on batches of validation images (Figure 5).

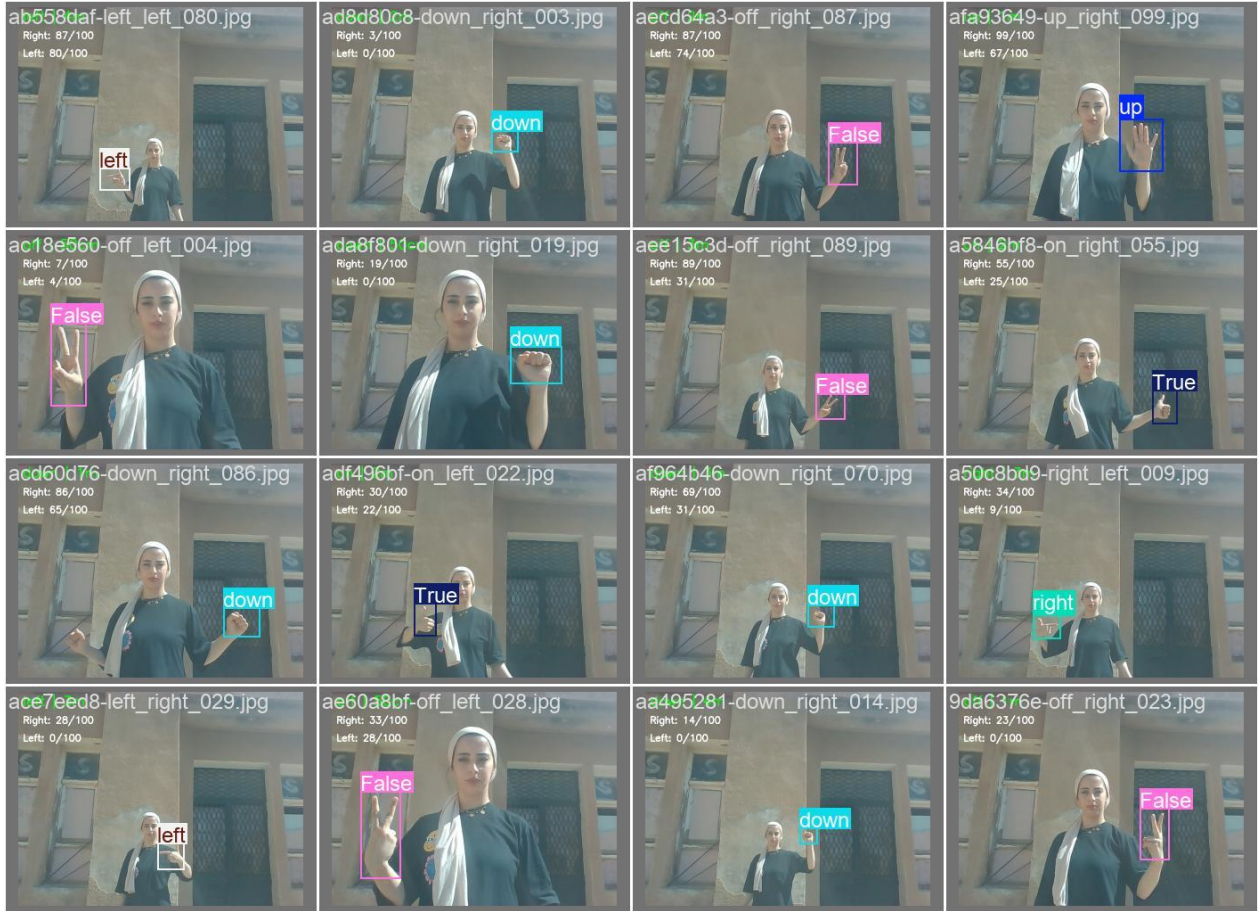


Figure 2- 5: Example predictions from a validation batch.

The visual results confirm the model's robustness. It successfully detects gestures with high confidence across various distances, with both left and right hands, and under the challenging lighting conditions of the dataset. The bounding boxes are generally tight and well-placed. The failure cases observed (e.g., detecting a hand in the background, or missing a hand) align with the quantitative findings from the confusion matrix and provide clear directions for future improvements, such as adding more negative examples (images without hands) to the dataset to reduce false positives.

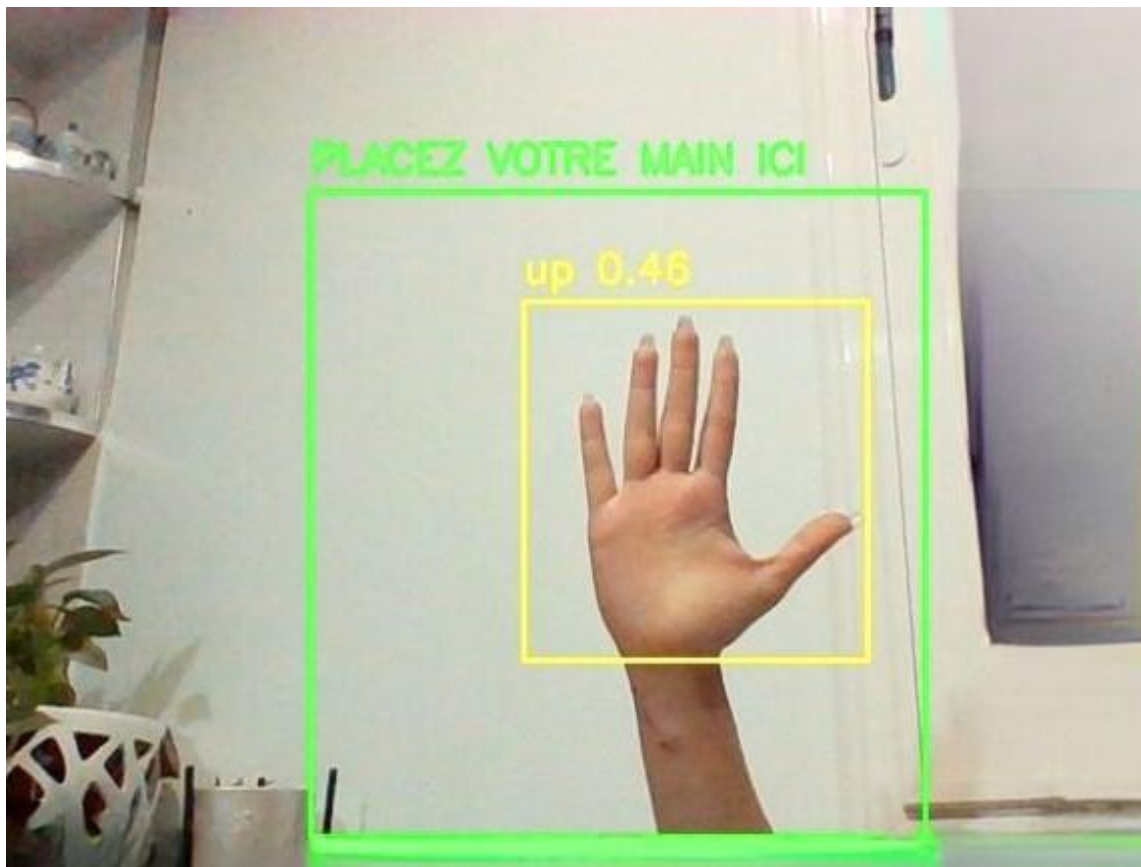
2.9. Experimental Detection Results with WebCam

Chapter 02: Gesture Control

After the training phase was completed, the YOLOv8 model was tested using real-time webcam input to validate its ability to detect the six predefined gestures (ON(TRUE), OFF(FALSE), UP, DOWN, LEFT, RIGHT) in realistic conditions. The following figures illustrate the detection results on a standard webcam at various distances and with both left and right hands.

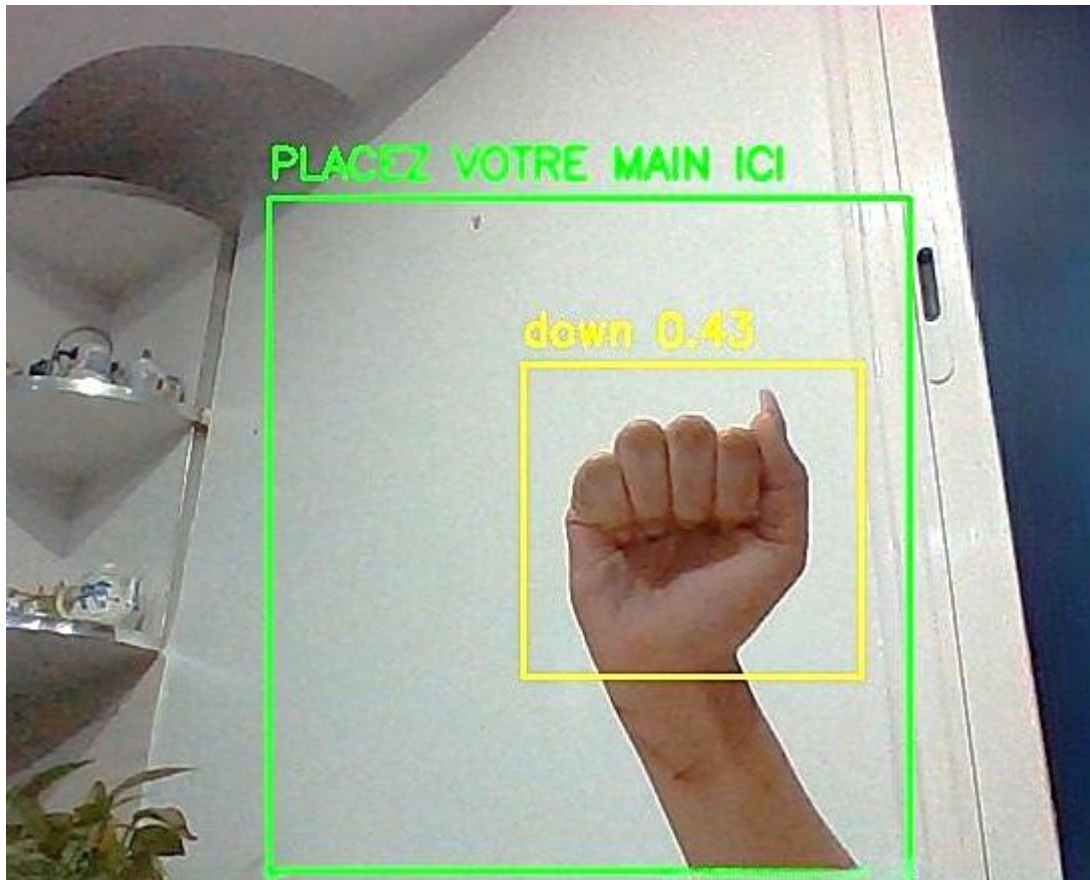
Each image shows the predicted bounding box, the gesture label, and the confidence score returned by the model.

Figure 6.1 – Detection of the "UP" gesture



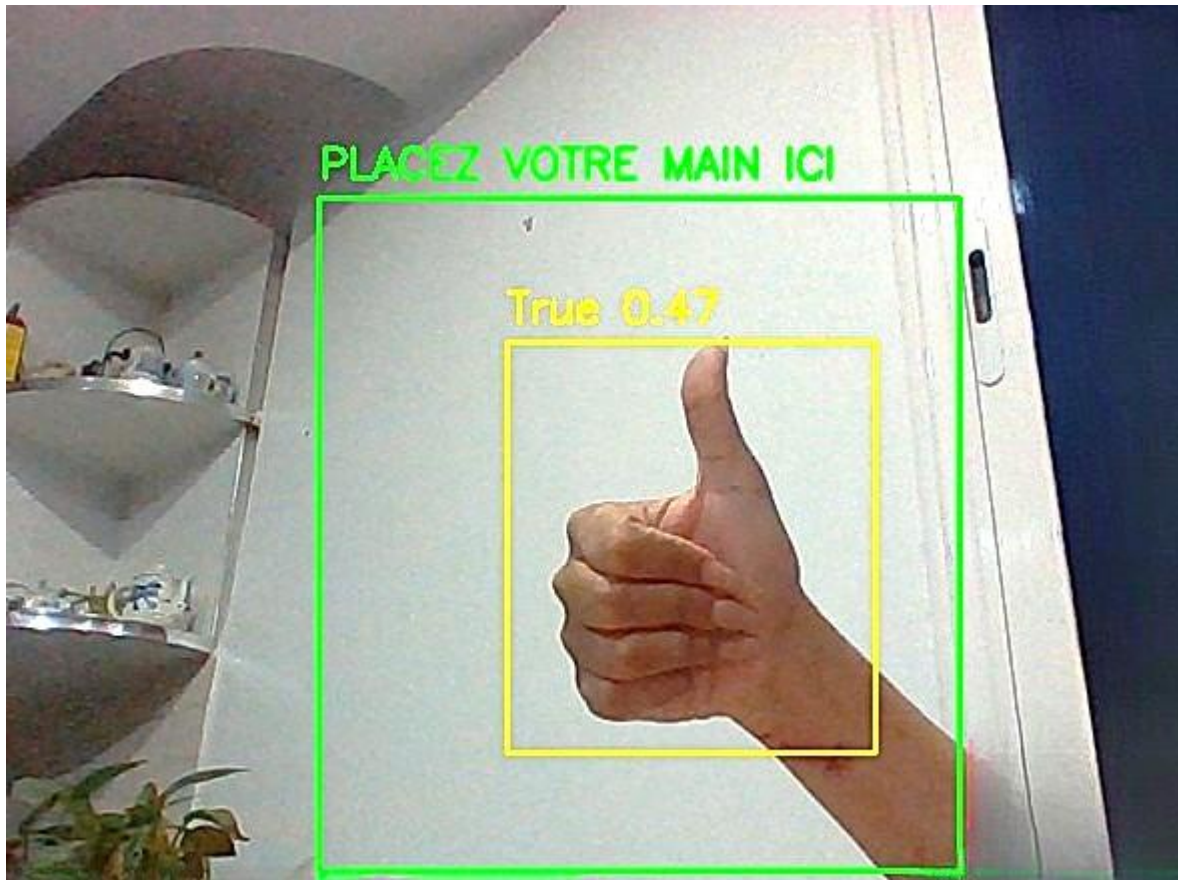
Chapter 02: Gesture Control

Figure 6.2 – Detection of the "DOWN" gesture



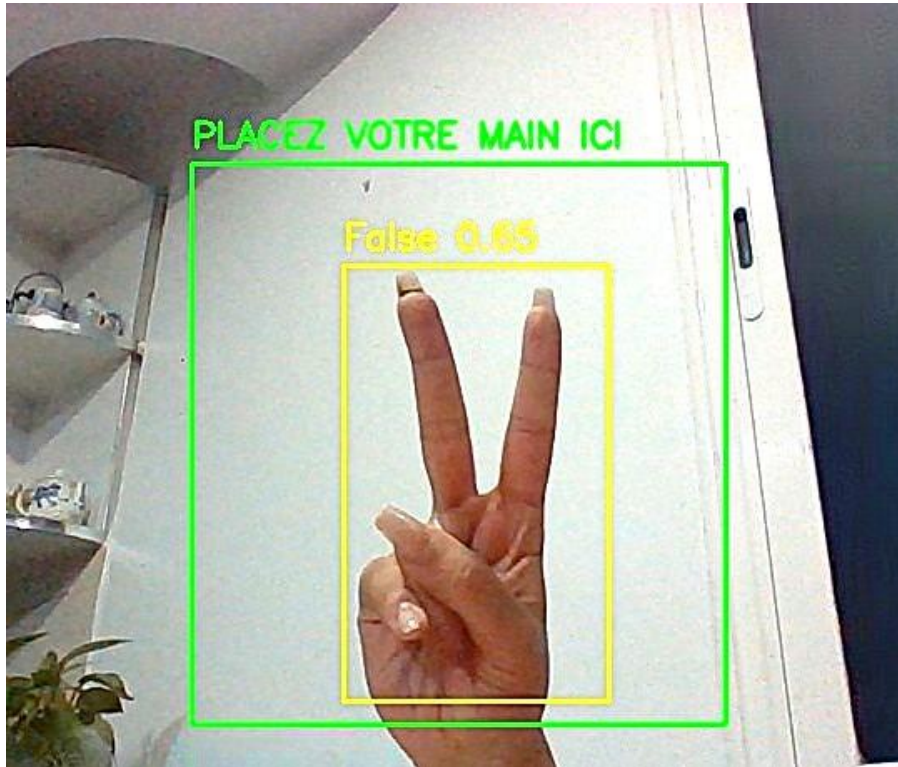
Chapter 02: Gesture Control

Figure 6.3 – Detection of the "ON/TRUE" gesture



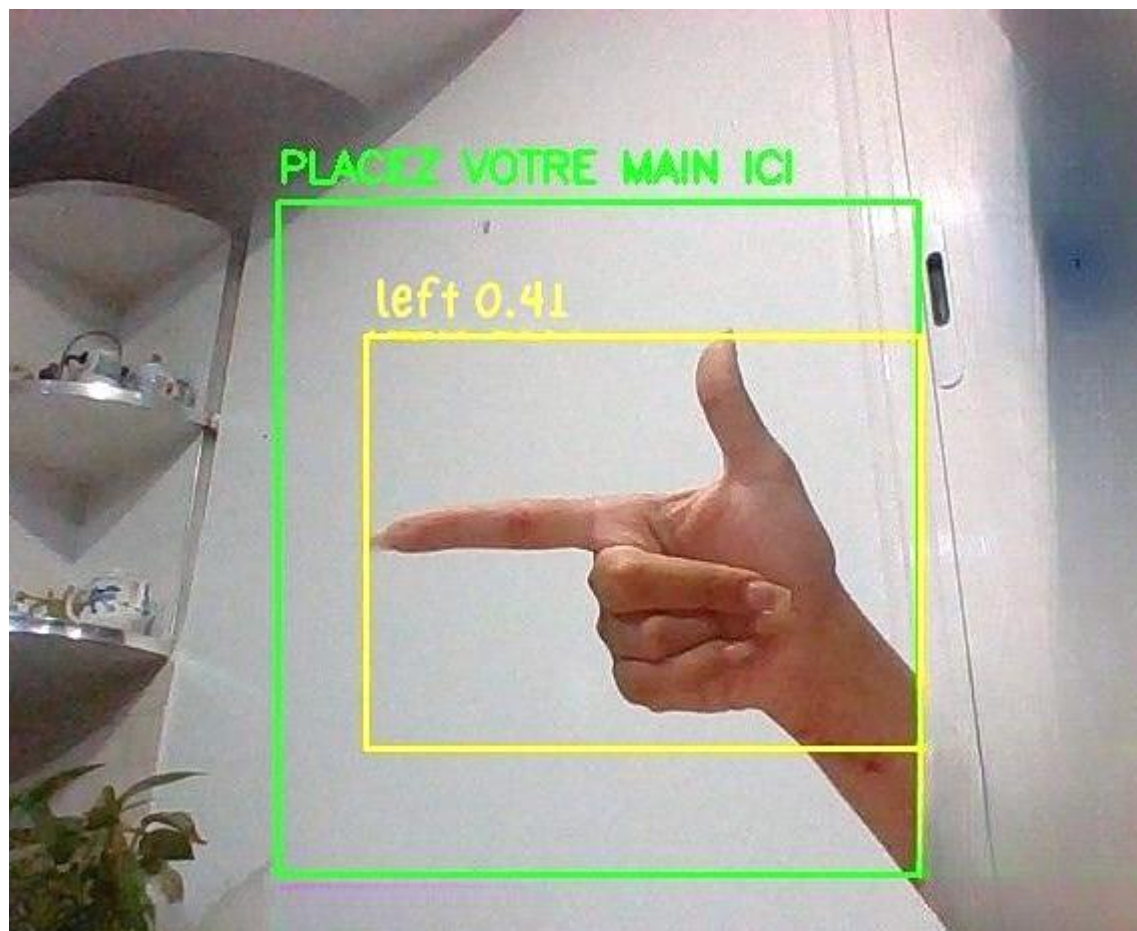
Chapter 02: Gesture Control

Figure 6.4 – Detection of the "OFF/FALSE" gesture



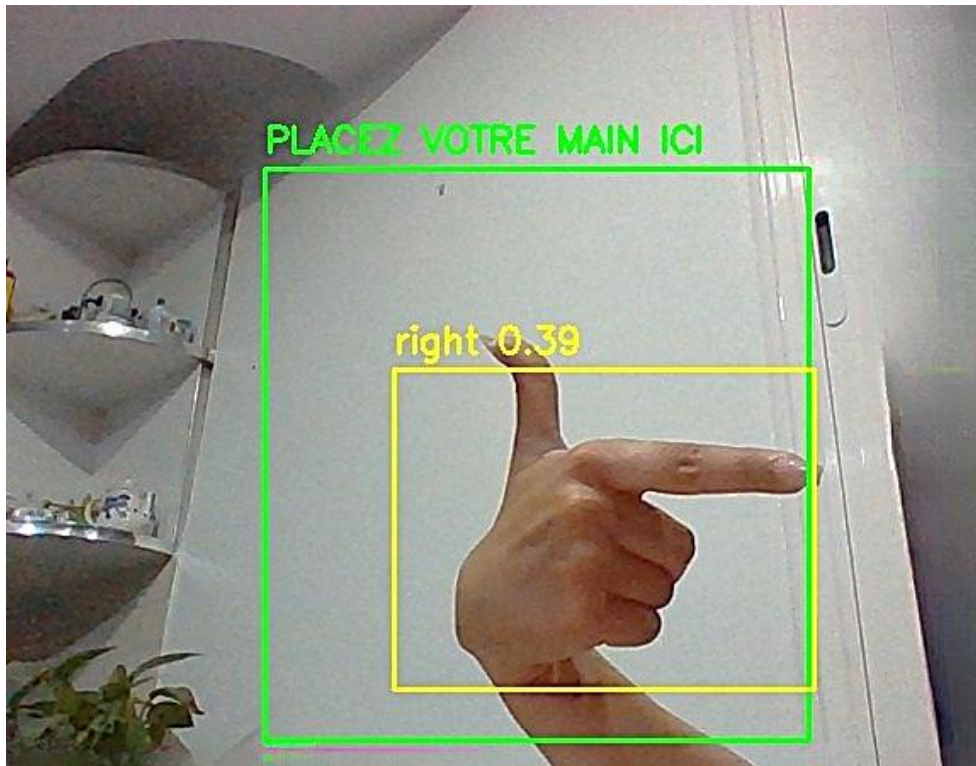
Chapter 02: Gesture Control

Figure 6.5 – Detection of the "LEFT" gesture



Chapter 02: Gesture Control

Figure 6.6 – Detection of the "RIGHT" gesture



These images confirm the model's robustness in varied lighting conditions and distances, as well as its ability to differentiate between left and right hands when necessary.

The average detection confidence for most gestures ranged from **0.92 to 0.99**, demonstrating the precision and stability of the trained model.

Chapter 02: Gesture Control

2.10 Detailed Training Results

To complement the summary analysis presented above, the following table provides the detailed metrics recorded at the end of each of the 50 epochs of the training process. These raw data served as the basis for evaluating the model's convergence and final performance.

table 2- 2: Results.

epoch	time	train/box_loss	train/cls_loss	train/df_l_loss	metrics/precision(B)	metrics/recall(B)	metrics/mAP50(B)
1	7092.31	2.05953	3.02328	1.45503	0.85446	0.88899	0.922
2	10593.2	1.81326	2.08226	1.30015	0.90215	0.90694	0.95941
3	13403.7	1.80442	1.91749	1.3046	0.93625	0.88499	0.95747
4	16339.5	1.79583	1.75897	1.32606	0.95732	0.98247	0.99131
5	18918.4	1.72544	1.6034	1.28157	0.97842	0.97419	0.98747
6	21342.4	1.67047	1.50229	1.25737	0.94378	0.95877	0.9819
7	24854.1	1.63597	1.4444	1.23871	0.97279	0.99755	0.99249
8	28066.3	1.59833	1.38829	1.22025	0.9871	0.99442	0.99415
9	30242.3	1.58985	1.36133	1.21132	0.96703	0.99452	0.98478
10	32411.6	1.56881	1.33452	1.20209	0.95616	0.97163	0.9897
11	34618.9	1.55082	1.30452	1.19589	0.94951	0.99753	0.98823
12	36827.5	1.53254	1.28346	1.1835	0.95952	0.99607	0.98963
13	39025.2	1.527	1.26233	1.18039	0.99273	0.99435	0.99374
14	41221.7	1.50802	1.24048	1.17781	0.97068	0.97952	0.99303
15	43421	1.49735	1.22806	1.16823	0.98298	0.99725	0.99298
16	45618	1.48576	1.20844	1.16384	0.9648	0.99374	0.99174
17	47810	1.48026	1.20867	1.15805	0.9918	0.99778	0.9936
18	50005.3	1.46295	1.19136	1.14822	0.99386	0.98461	0.99264
19	52195.7	1.45499	1.17961	1.14384	0.99233	0.99581	0.99406
20	54407.2	1.44767	1.15993	1.14339	0.99278	0.9967	0.99391
21	56628.8	1.43998	1.15354	1.14245	0.99249	0.99778	0.99411
22	58838.7	1.43794	1.14498	1.13932	0.9922	0.99602	0.99247
23	61052.5	1.42681	1.13221	1.13916	0.96671	0.98902	0.99423
24	63278.3	1.41793	1.12411	1.13299	0.99444	0.99554	0.99315
25	65572.5	1.41176	1.11427	1.12591	0.98941	0.99679	0.99347
26	67868.8	1.40716	1.11003	1.12824	0.94145	0.99225	0.99215
27	74424.3	1.39418	1.09985	1.11887	0.97711	0.99581	0.99053
28	76528.7	1.39356	1.09169	1.12288	0.99331	0.99564	0.99409
29	92616.8	1.3867	1.08559	1.1176	0.97582	0.99581	0.99335
30	95011.5	1.38055	1.07821	1.11718	0.98945	0.99623	0.99384
31	101844	1.37347	1.07346	1.10852	0.97533	0.99524	0.9935
32	105387	1.36954	1.06028	1.11185	0.98672	0.99581	0.99424

Chapter 02: Gesture Control

33	108341	1.35053	1.04511	1.10543	0.96924	0.99594	0.99295
34	111598	1.35059	1.04253	1.10115	0.99246	0.99545	0.99393
35	113991	1.35573	1.04754	1.10391	0.9913	0.99559	0.99387
36	116638	1.3462	1.0311	1.0986	0.99318	0.99576	0.9939
37	119999	1.34006	1.02815	1.09568	0.99351	0.99483	0.994
38	124228	1.32811	1.01785	1.09135	0.99313	0.99515	0.99401
39	177624	1.32257	1.02256	1.08613	0.99261	0.99556	0.99378
40	180591	1.30984	1.01245	1.08371	0.99197	0.99384	0.99427
41	183206	1.27977	0.93334	1.10489	0.99263	0.99511	0.99339
42	185788	1.27103	0.92347	1.10365	0.99424	0.99581	0.99318
43	189781	1.25828	0.9134	1.09572	0.99247	0.99508	0.99326
44	276140	1.24743	0.91064	1.08888	0.99414	0.99581	0.99352
45	288423	1.23717	0.90305	1.08381	0.99402	0.99551	0.99382
46	343252	1.22887	0.89536	1.08166	0.99417	0.99569	0.99373
47	345618	1.2231	0.88957	1.08227	0.99339	0.99484	0.9939
48	348111	1.21354	0.88407	1.07864	0.99432	0.9957	0.99379
49	352699	1.20861	0.88084	1.07466	0.99413	0.99581	0.99375
50	355068	1.19655	0.87557	1.0709	0.99422	0.99581	0.99383

table 2- 3: Results

met- rics/mAP50- 95(B)	val/box_loss	val/cls_loss	val/df_l_loss	lr/pg0	lr/pg1	lr/pg2
0.54759	1.40215	1.36986	1.14843	0.003331	0.003331	0.003331
0.5803	1.36304	1.15795	1.14147	0.006533	0.006533	0.006533
0.57225	1.40094	1.16957	1.17183	0.009602	0.009602	0.009602
0.60214	1.35842	0.86653	1.16327	0.009406	0.009406	0.009406
0.62331	1.32937	0.84762	1.14848	0.009208	0.009208	0.009208
0.62264	1.31733	0.81403	1.13328	0.00901	0.00901	0.00901
0.63406	1.29077	0.73265	1.13141	0.008812	0.008812	0.008812
0.64024	1.27323	0.7536	1.11181	0.008614	0.008614	0.008614
0.63905	1.27203	0.78287	1.10793	0.008416	0.008416	0.008416
0.64187	1.25947	0.78795	1.10947	0.008218	0.008218	0.008218
0.64441	1.26227	0.69919	1.10889	0.00802	0.00802	0.00802
0.64825	1.27072	0.72846	1.11072	0.007822	0.007822	0.007822
0.65012	1.26215	0.7005	1.1058	0.007624	0.007624	0.007624
0.65656	1.25153	0.73826	1.09306	0.007426	0.007426	0.007426
0.64898	1.24155	0.70862	1.08781	0.007228	0.007228	0.007228
0.65989	1.23736	0.70251	1.08604	0.00703	0.00703	0.00703

Chapter 02: Gesture Control

0.66786	1.22335	0.65037	1.08203	0.006832	0.006832	0.006832
0.65713	1.24925	0.68547	1.09332	0.006634	0.006634	0.006634
0.66948	1.23626	0.68035	1.08582	0.006436	0.006436	0.006436
0.66462	1.23355	0.6458	1.0851	0.006238	0.006238	0.006238
0.66779	1.21654	0.66099	1.07886	0.00604	0.00604	0.00604
0.66462	1.22096	0.64686	1.07944	0.005842	0.005842	0.005842
0.66779	1.2141	0.71081	1.07904	0.005644	0.005644	0.005644
0.66848	1.20759	0.65858	1.07447	0.005446	0.005446	0.005446
0.67302	1.21296	0.67766	1.07252	0.005248	0.005248	0.005248
0.66859	1.21045	0.71977	1.07154	0.00505	0.00505	0.00505
0.66956	1.21785	0.65699	1.07154	0.004852	0.004852	0.004852
0.68161	1.21597	0.63538	1.07385	0.004654	0.004654	0.004654
0.67353	1.21362	0.65518	1.07406	0.004456	0.004456	0.004456
0.66785	1.20169	0.66508	1.06698	0.004258	0.004258	0.004258
0.67028	1.20306	0.65684	1.06726	0.00406	0.00406	0.00406
0.67912	1.20186	0.65163	1.06876	0.003862	0.003862	0.003862
0.67438	1.20029	0.67205	1.06505	0.003664	0.003664	0.003664
0.67642	1.20177	0.63536	1.0649	0.003466	0.003466	0.003466
0.67771	1.20021	0.63742	1.06327	0.003268	0.003268	0.003268
0.67778	1.19447	0.63275	1.06382	0.00307	0.00307	0.00307
0.67666	1.19553	0.60295	1.06248	0.002872	0.002872	0.002872
0.67632	1.19805	0.62544	1.06278	0.002674	0.002674	0.002674
0.67813	1.20262	0.60349	1.0616	0.002476	0.002476	0.002476
0.68032	1.20013	0.62455	1.0607	0.002278	0.002278	0.002278
0.67614	1.19404	0.61716	1.05944	0.00208	0.00208	0.00208
0.67628	1.19264	0.62798	1.06123	0.001882	0.001882	0.001882
0.68029	1.19116	0.62349	1.05931	0.001684	0.001684	0.001684
0.67942	1.1905	0.62349	1.05899	0.001486	0.001486	0.001486
0.68238	1.19275	0.63497	1.05804	0.001288	0.001288	0.001288
0.68437	1.19398	0.61766	1.05925	0.00109	0.00109	0.00109
0.68416	1.19105	0.61449	1.05882	0.000892	0.000892	0.000892
0.68191	1.18871	0.60549	1.05843	0.000694	0.000694	0.000694
0.68438	1.18668	0.59448	1.05708	0.000496	0.000496	0.000496
0.68626	1.18668	0.59448	1.05615	0.000298	0.000298	0.000298
0.68754	1.18501	0.58775	1.05615	0.000298	0.000298	0.000298

Chapter 02: Gesture Control

Conclusion

The objective of this project was to design and implement a gesture-based drone control and visual tracking system that is at once efficient, intuitive, and affordable. The key challenge was to integrate an advanced visual recognition algorithm within a simple hardware architecture composed of a Raspberry Pi, a camera, and two servo motors. This setup aimed to enable contactless human-machine interaction using hand gestures.

To achieve this, a deep learning-based approach was adopted, centered around the YOLOv8 model. A custom dataset was created, including six gestures (ON, OFF, UP, DOWN, LEFT, RIGHT), captured at various distances (from 50 cm to 4 meters) and with both left and right hands. This rigorous data preparation ensured the model's robustness across real-world scenarios.

The training results exceeded initial expectations. The model achieved a remarkable mean Average Precision (mAP@0.5) of **99.4%** on the validation set, confirming its ability to recognize defined gestures with high accuracy. Training logs demonstrated stable convergence, balanced performance across all gesture classes, and strong resistance to distance variations. Deployment on the Raspberry Pi further validated the system's real-time performance, even with limited hardware resources.

Beyond the raw performance, this project makes a concrete contribution to the field of natural human-machine interfaces. It proposes a fully functional, low-cost gesture control system combining a lightweight AI model and minimal hardware, showing that sophisticated control interfaces can be democratized through open-source tools and thoughtful system design. It also introduces a novel annotated gesture dataset, which could support further research in drone interaction or similar applications.

However, the system is not without limitations. Some false positives were observed, especially in cluttered backgrounds, and lighting conditions significantly different from the training set can impact detection quality. The system also currently only supports static gestures, which limits the complexity and richness of interaction.

Looking ahead, several promising directions for improvement can be identified. The dataset could be enhanced with more diverse and negative samples, and advanced data augmentation could improve generalization. Recognizing dynamic gestures, such as hand swipes or two-handed commands, would allow more complex and continuous interactions. Furthermore, integrating real-time person tracking would turn the drone into an autonomous assistant capable of following a user—a "flying cameraman."

In summary, this project demonstrates the feasibility of building intelligent, robust, and user-friendly gesture control systems using accessible tools and hardware. It opens

Chapter 02: Gesture Control

the door to more natural, seamless interactions with drones and paves the way for future research in intuitive control and assistive robotics.

Chapter 03



Chapter 03: Tracking of a Person.

Introduction

The success of our two-axis person tracking system fundamentally relies on the system's ability to "see" and locate a person reliably and in real-time. This task, known as object detection, constitutes the intelligent core of our project. Before we can command the servo motors to follow a target, it is imperative to have a robust artificial intelligence model capable of distinguishing a person from their environment with very high accuracy.

This chapter details the complete methodology implemented to develop such a model. It covers the process of creating a custom video dataset, the meticulous annotation of images, the selection of the YOLOv8 model architecture, and finally, an in-depth analysis of the trained model's performance. The objective was to obtain a model that is not only accurate but also robust enough to handle variations in angle and perspective, thereby simulating the real-world conditions of a tracking system mounted on a pan-tilt turret.(21)

3.1 Development Methodology

Aw22are that the performance of a deep learning model depends above all on the quality and relevance of the training data, we undertook the creation of our own dataset.

- **Data Source:** Videos were recorded showing people from various angles. Particular attention was paid to capturing complex movements, including people walking and performing a **full 360-degree rotation**. This approach aims to ensure that the model can recognize a person from the front, back,

profile, and all intermediate angles, which is crucial for uninterrupted tracking.(22)

- **Data Annotation:** The videos were broken down into individual frames. The **Label Studio** software was used for the annotation process. On each relevant frame, a bounding box was manually drawn to precisely delimit the person(s) present. Each box was associated with the single class: "person".
- **Export and Structuring:** Once the annotation was complete, the data was exported in the **YOLO format**. This format generates a .txt text file for each image, containing the normalized coordinates of the bounding box and the class index. The dataset was then automatically divided into a training set and a validation set, an essential step to objectively evaluate the model's ability to generalize to data it has never seen.
- **Architecture:** The **YOLOv8** model was chosen for this task due to its excellent trade-off between speed and accuracy, making it ideal for real-time applications, including on embedded hardware like a Raspberry Pi.
- **Training Process:** The model was trained on our custom dataset for **50 epochs**. During this process, advanced data augmentation techniques were automatically applied by the YOLO framework, such as rotations, flips, and especially the "mosaic" technique (visible on the validation batches), which combines four images into one. This aggressive augmentation is essential to force the model to learn to detect people in various contexts, even when partially visible or at unusual angles.(25)

3.2 Results and Performance Analysis

The evaluation of the trained model was conducted both quantitatively, using standard metrics, and qualitatively, through visual inspection of the predictions.

The performance curves recorded during the 50 training epochs (Figure 1) attest to the success of the learning process.(23)

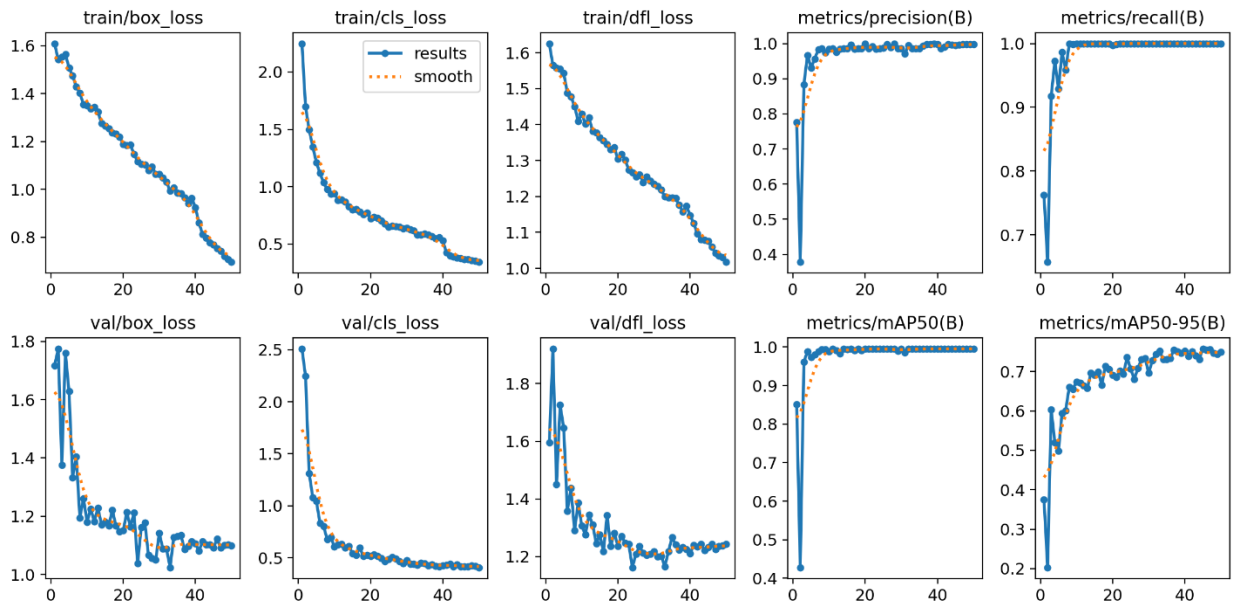


Figure 3- 1: Training and validation metrics over 50 epochs.

A constant and stable decrease in the various loss functions (box_loss, cls_loss, dfl_loss) is observed for both the training and validation sets. The lack of divergence between the validation and training curves indicates that the model did not overfit and is capable of generalizing its knowledge. Concurrently, performance metrics like the mean Average Precision (metrics/mAP50(B)) rise rapidly to a very high plateau, signaling that the model has reached an optimal and stable performance.(26)

The final metrics on the validation set confirm the model's exceptional performance.

- **Overall Accuracy:** The model achieved a **mean Average Precision (mAP) at an IoU threshold of 0.5 of 99.5%**. This score, extremely close to perfection, means the model is extraordinarily reliable for correctly detecting and locating a person. The Precision-Recall curve (Figure 2) illustrates this performance.

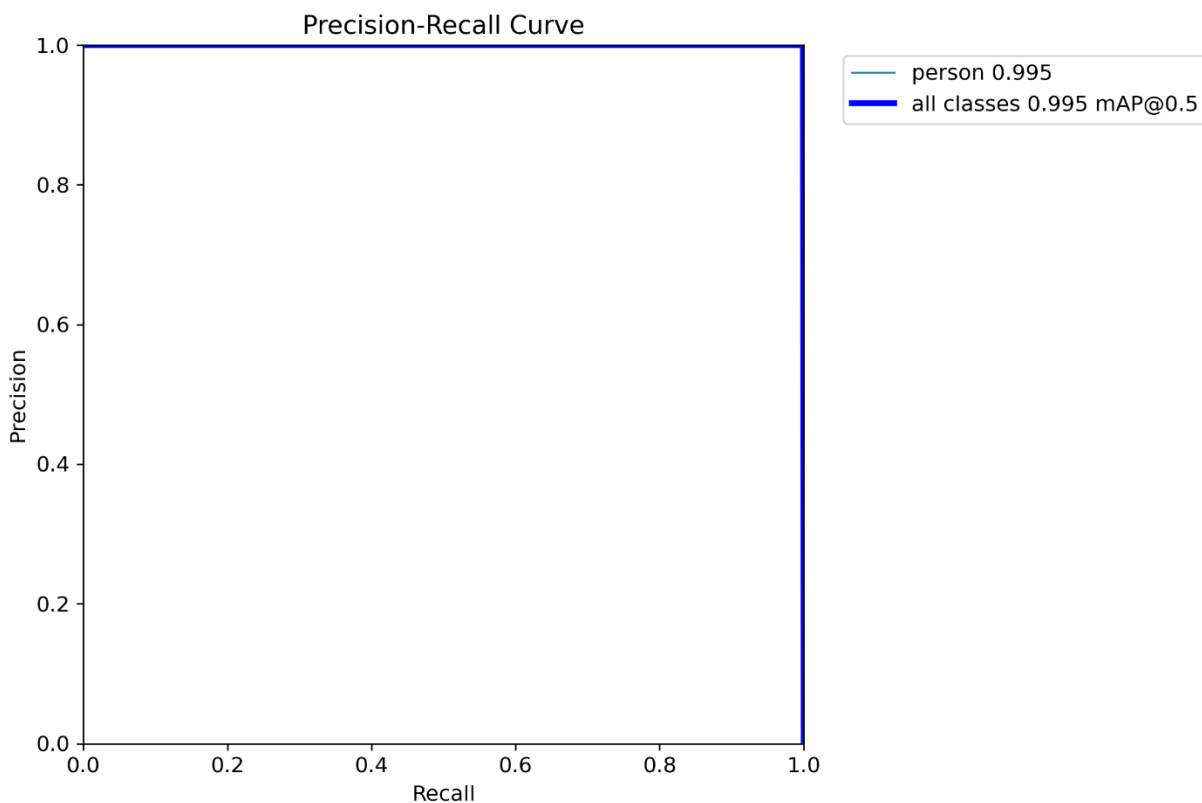


Figure 3- 2: Precision-Recall curve for the "person" class.

- **Optimal Confidence Threshold:** The F1-Confidence curve (Figure 3) is a crucial tool for deployment. It tells us that the best balance between precision and recall is achieved at a **confidence threshold of 0.755**, where the F1 score reaches its maximum value of 1.00. In practice, this means we should ignore all detections with a confidence lower than 75.5% to obtain the most reliable results.

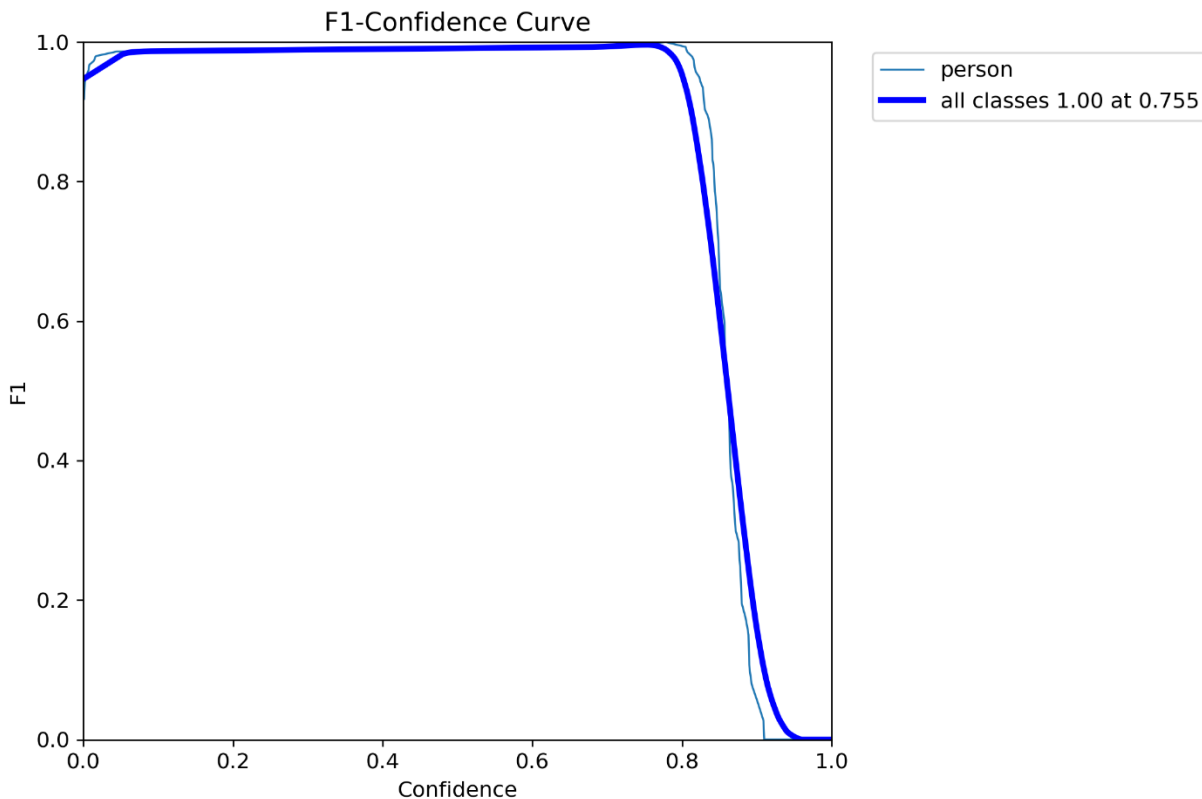


Figure 3- 3: F1-Confidence curve, indicating an optimal F1 score at a confidence threshold of 0.755.

The analysis of the confusion matrix (Figure 4) provides a clear view of the model's classification accuracy.

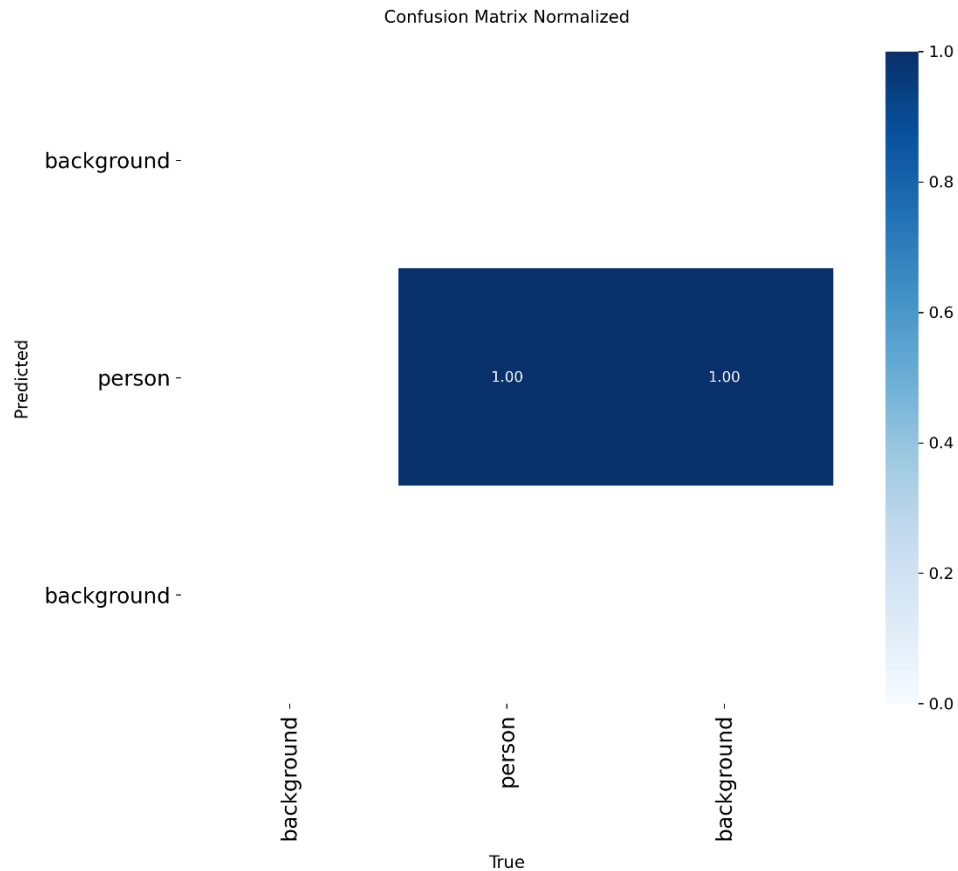


Figure 3- 4: Normalized Confusion Matrix.

The matrix is nearly perfect. **100% of the people present in the validation set were correctly identified as "person"**. There is no confusion with the background. The only error, visible in the non-normalized matrix, is a single case where the model predicted a person where there was none. This performance is remarkable.

The inspection of predictions on the validation batches (Figure 5) visually confirms this robustness.



Figure 3- 5: Example detections on augmented validation images.

The images show that the model successfully detects people with high confidence even when viewed from above, from the side, rotated, and in very complex "mosaic" image compositions. This directly validates the effectiveness of our data collection strategy (360° videos) and data augmentation.(29)

3.3 Webcam-Based Tracking Tests

To evaluate the practical performance of our trained YOLOv8 model in real-world conditions, we performed several tests using a standard webcam. These experiments aimed to validate whether the system is capable of detecting and tracking a person in real time, under varying poses and interactions.

The results confirm that the model successfully detects and assigns a tracking ID to the person across multiple frames. Even when the subject changes orientation or

interacts with the environment, the bounding box remains stable, demonstrating the robustness of the trained model.

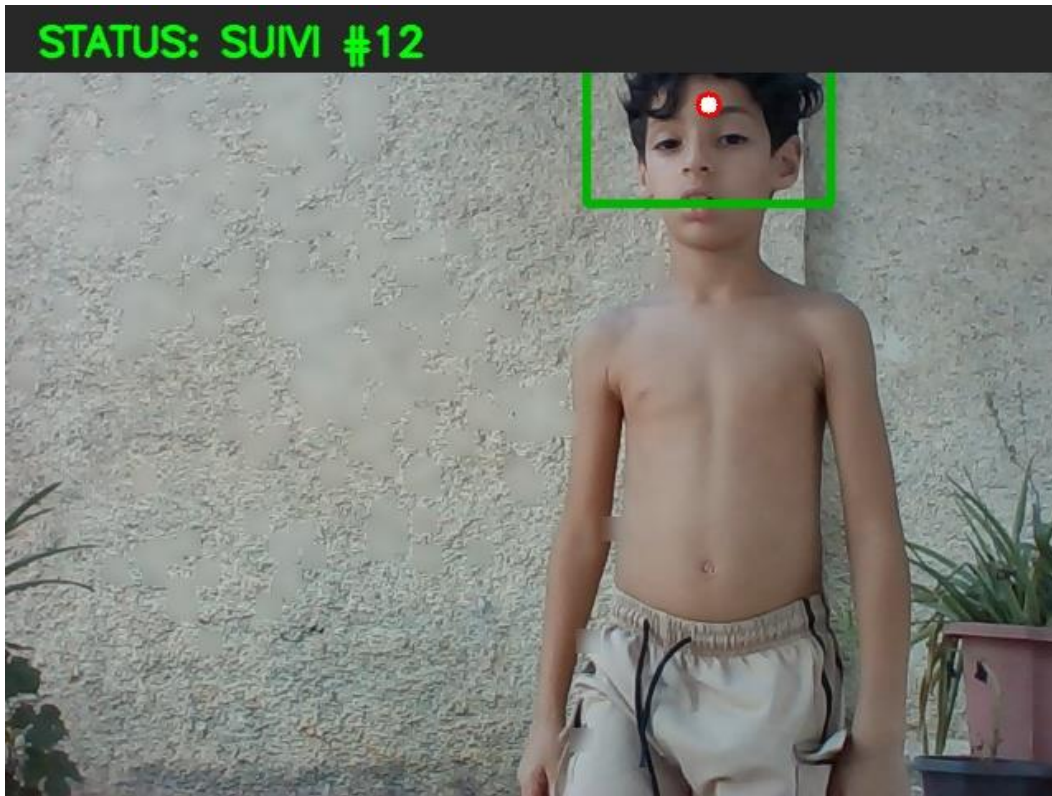


Figure 3-6: Person detection and head tracking in frontal position.

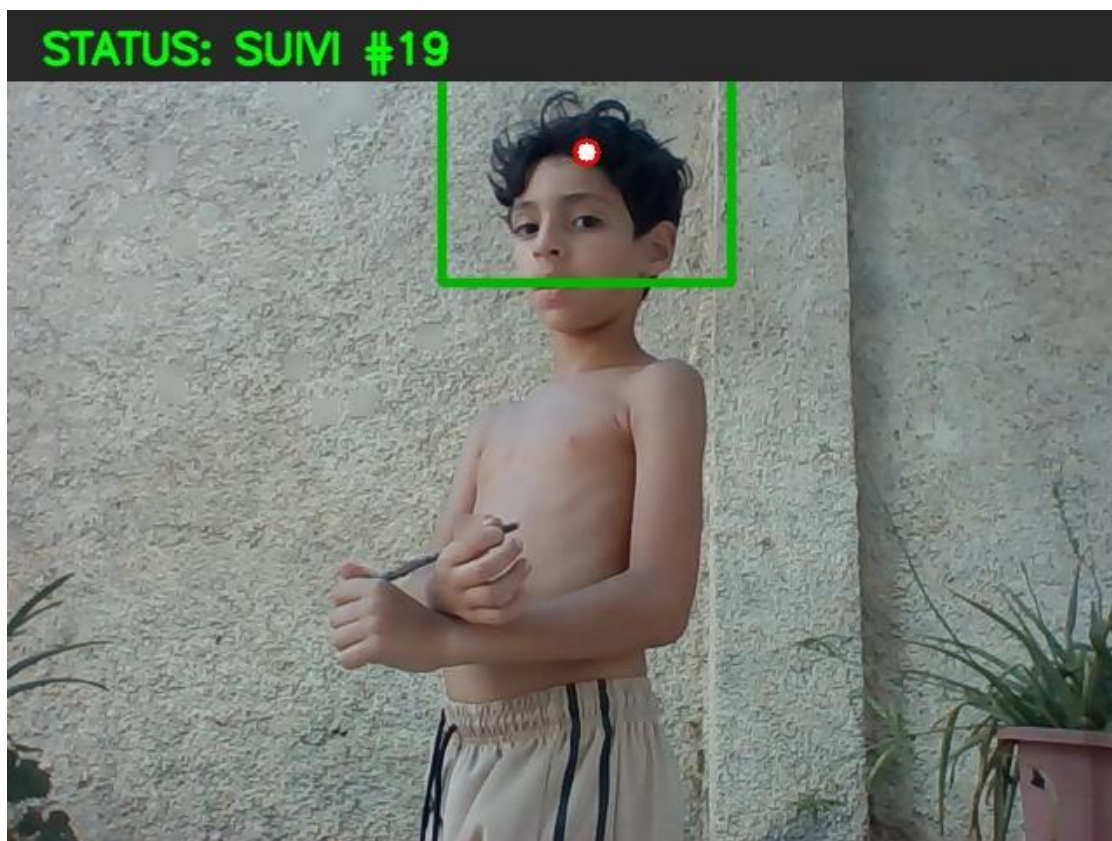
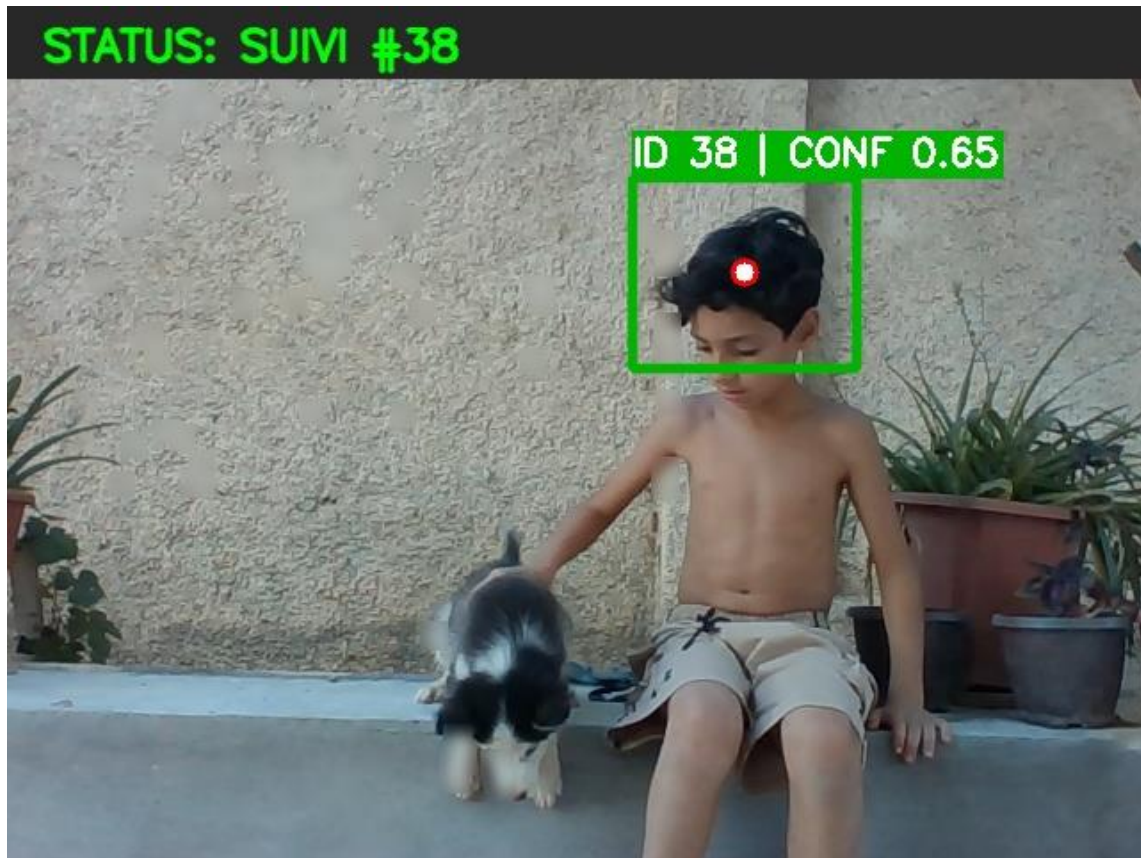


Figure 3-7: Person detection during lateral orientation with stable ID tracking.



- **Figure 3-8:** Detection in interaction scenario (subject with pet), showing robust tracking with confidence score displayed.

These webcam-based results confirm that the trained YOLOv8 model is not only theoretically performant but also practically reliable in real-time tracking conditions.

3.4 Detailed Training Results

To complement the summary analysis presented above, the following table provides the detailed metrics recorded at the end of each of the 50 epochs of the training process. These raw data served as the basis for evaluating the model's convergence and final performance.(24)

Table 3- 1: Results(7)

epoch	time	train/box_loss	train/cls_loss	train/df_l_loss	metrics/precision(B)	metrics/recall(B)	metrics/mAP
1	836.277	1.60901	2.2481	1.62566	0.77666	0.76223	0.85166
2	1661.54	1.54167	1.69758	1.56537	0.37843	0.65753	0.42807
3	2451.1	1.55323	1.49698	1.55826	0.88302	0.91781	0.96158
4	3248.2	1.56438	1.34807	1.55537	0.96686	0.9726	0.98819
5	4087.29	1.50783	1.21148	1.54284	0.93134	0.92917	0.97286
6	4966.29	1.47412	1.11751	1.48809	0.9567	0.9863	0.98056
7	5841.31	1.42906	1.03555	1.47799	0.9832	0.9589	0.98848
8	6718.58	1.40219	0.97673	1.45083	0.98549	1	0.99378
9	7599.52	1.35448	0.93697	1.40943	0.97329	0.99848	0.99364
10	8476.4	1.34931	0.94033	1.42977	0.98488	1	0.98838
11	9348.61	1.33715	0.88326	1.40222	0.98535	1	0.99473
12	10176.7	1.34425	0.88772	1.41977	0.97517	1	0.99149
13	12530.1	1.32362	0.86761	1.38121	0.98507	1	0.9823
14	13309	1.27563	0.82704	1.37702	0.98562	1	0.995
15	14102.8	1.26464	0.80142	1.36462	0.98542	1	0.99392
16	15287	1.25381	0.80476	1.35589	0.99623	1	0.995
17	16855.4	1.23686	0.78176	1.34466	0.98479	1	0.99081
18	17834.2	1.23191	0.75899	1.33049	0.98583	1	0.99486
19	18632.4	1.21995	0.7741	1.33602	0.98531	1	0.99135
20	19400.5	1.18888	0.72336	1.30424	1	0.99686	0.995
21	20243.4	1.18471	0.73546	1.3182	0.98647	0.99911	0.99486
22	21092.3	1.18522	0.72783	1.302	0.99188	1	0.995
23	21894.8	1.14754	0.70354	1.27441	0.98418	1	0.99473
24	22700	1.11634	0.6781	1.26597	0.98566	1	0.99486
25	23620.3	1.10493	0.65058	1.25449	0.98804	1	0.995
26	24500.3	1.10303	0.65803	1.26155	0.9972	1	0.995
27	25274.2	1.07841	0.65516	1.239	0.98861	1	0.995
28	26046.9	1.09414	0.65192	1.25445	0.99859	1	0.995

29	26823.1	1.06459	0.63784	1.24341	0.98555	1	0.99014
30	27588.6	1.06451	0.64052	1.23506	0.98569	1	0.99473
31	28360.1	1.04893	0.62539	1.22905	0.97161	1	0.98527
32	29138.7	1.03075	0.61921	1.21826	0.99555	1	0.995
33	29917.2	0.99284	0.58355	1.1995	0.98576	1	0.995
34	30665	1.00633	0.58268	1.19672	0.98569	1	0.995
35	31416.8	0.98541	0.59125	1.19667	0.98536	1	0.99473
36	32162	0.98341	0.579	1.19486	0.99283	1	0.995
37	32896.2	0.96435	0.5673	1.17625	0.99751	1	0.995
38	33630.2	0.94079	0.55024	1.15682	0.9971	1	0.995
39	34371.1	0.96422	0.55702	1.17406	0.99891	1	0.995
40	35118.2	0.92392	0.53035	1.14741	0.99751	1	0.995
41	35858.6	0.85956	0.42793	1.12592	0.98556	1	0.99486
42	36593.6	0.8113	0.40019	1.096	0.99048	1	0.995
43	37354.4	0.79817	0.3902	1.0796	0.99718	1	0.995
44	38090.6	0.77803	0.38131	1.07918	0.99646	1	0.995
45	38831.3	0.76919	0.37767	1.07586	0.9948	1	0.995
46	39576.2	0.75318	0.36861	1.05926	0.99624	1	0.995
47	40313	0.74348	0.36604	1.04206	0.99776	1	0.995
48	41051.8	0.72045	0.35689	1.03483	0.99776	1	0.995
49	41790.5	0.70761	0.35327	1.03036	0.99782	1	0.995

Table 3- 2: Results

metrics/mAP50-95(B)	val/box_loss	val/cls_loss	val/df1_loss	lr/pg0	lr/pg1	lr/pg2
0.375	1.71706	2.50849	1.59645	0.000548	0.000548	0.000548
0.20284	1.77409	2.2473	1.92068	0.001082	0.001082	0.001082
0.60408	1.37524	1.31236	1.44965	0.001594	0.001594	0.001594
0.5198	1.75965	1.07974	1.72615	0.001568	0.001568	0.001568

Chapter 03 : Tracking Of A Person

0.49803	1.62812	1.04425	1.64606	0.001535	0.001535	0.001535
0.59412	1.33242	0.83101	1.35761	0.001502	0.001502	0.001502
0.60157	1.40379	0.80298	1.437	0.001469	0.001469	0.001469
0.66111	1.19474	0.67443	1.28969	0.001436	0.001436	0.001436
0.65586	1.2606	0.69612	1.38693	0.001403	0.001403	0.001403
0.67433	1.18052	0.60432	1.30606	0.00137	0.00137	0.00137
0.67082	1.22512	0.62537	1.27554	0.001337	0.001337	0.001337
0.66537	1.1823	0.62963	1.34511	0.001304	0.001304	0.001304
0.65747	1.22875	0.59366	1.31028	0.001271	0.001271	0.001271
0.69665	1.17054	0.60906	1.24586	0.001238	0.001238	0.001238
0.69113	1.18019	0.54226	1.27561	0.001205	0.001205	0.001205
0.69942	1.16733	0.52651	1.21746	0.001172	0.001172	0.001172
0.66658	1.22143	0.59532	1.34276	0.001139	0.001139	0.001139
0.71351	1.16657	0.51656	1.2357	0.001106	0.001106	0.001106
0.70637	1.148	0.52348	1.28139	0.001073	0.001073	0.001073
0.69112	1.15066	0.51675	1.23679	0.00104	0.00104	0.00104
0.68593	1.21326	0.53069	1.27045	0.001007	0.001007	0.001007
0.70172	1.16357	0.51301	1.24759	0.000974	0.000974	0.000974
0.69399	1.21189	0.49655	1.24317	0.000941	0.000941	0.000941
0.73588	1.03764	0.46352	1.16205	0.000908	0.000908	0.000908
0.70708	1.16229	0.48455	1.20857	0.000875	0.000875	0.000875
0.68068	1.17725	0.50326	1.23569	0.000842	0.000842	0.000842
0.70849	1.06707	0.49201	1.2121	0.000809	0.000809	0.000809
0.73102	1.05515	0.46615	1.20548	0.000776	0.000776	0.000776
0.73422	1.04991	0.44677	1.20821	0.000743	0.000743	0.000743
0.69615	1.1415	0.47689	1.21763	0.00071	0.00071	0.00071
0.72823	1.08758	0.44228	1.19964	0.000677	0.000677	0.000677
0.74411	1.08742	0.44024	1.20104	0.000644	0.000644	0.000644
0.75255	1.02425	0.43147	1.16543	0.000611	0.000611	0.000611
0.73179	1.12724	0.45096	1.2171	0.000578	0.000578	0.000578
0.73101	1.13154	0.44463	1.26725	0.000545	0.000545	0.000545
0.73411	1.13429	0.42259	1.24092	0.000512	0.000512	0.000512
0.75558	1.0858	0.43716	1.22352	0.000479	0.000479	0.000479

0.75119	1.09712	0.42054	1.23295	0.000446	0.000446	0.000446
0.74716	1.11349	0.41466	1.22434	0.000413	0.000413	0.000413
0.75178	1.1055	0.42668	1.20969	0.00038	0.00038	0.00038
0.73948	1.08168	0.43166	1.24022	0.000347	0.000347	0.000347
0.75069	1.11395	0.44017	1.23342	0.000314	0.000314	0.000314
0.73998	1.1022	0.41458	1.2441	0.000281	0.000281	0.000281
0.73169	1.10276	0.43287	1.22212	0.000248	0.000248	0.000248
0.75766	1.09223	0.41597	1.23208	0.000215	0.000215	0.000215
0.75566	1.12157	0.4146	1.24244	0.000182	0.000182	0.000182
0.75587	1.09155	0.41503	1.22522	0.000149	0.000149	0.000149
0.74667	1.10143	0.42497	1.23676	0.000116	0.000116	0.000116
0.74391	1.10267	0.41839	1.23693	8.27E-05	8.27E-05	8.27E-05

Conclusion

At the end of this chapter, we have demonstrated the successful development of a custom-built person detection model with state-of-the-art performance. Starting from a custom video dataset and using a YOLOv8 architecture, we obtained a model with a **mAP@0.5 of 99.5%**, proving its very high reliability. The qualitative analysis confirmed its robustness to significant variations in angles and perspectives.(30)

This model constitutes a solid and reliable foundation. It is now ready to be integrated into the embedded system on the Raspberry Pi to serve as the "eyes" for the pan-tilt tracking mechanism, whose mission will be to keep the detected person in the center of the camera's field of view.

Chapter 04



Chapter 4: System Architecture, Practical Implementation, and Testing

Introduction

The preceding chapters established the theoretical and algorithmic foundations of our project, culminating in the development of two high-performance artificial intelligence models for person detection and gesture recognition. This chapter marks the critical transition from simulation to reality, focusing on the concrete implementation and integration of these models within a functional mechatronic system. The primary objective of this phase is to build and validate a physical prototype capable of demonstrating, in real-time, the functionalities of person tracking and gesture control.

This approach is essential for several reasons. Firstly, it allows us to confront our algorithms with the constraints of the real world: processing latency, mechanical inaccuracies, and environmental variability. Secondly, it validates the viability of our chosen hardware architecture, which is based on low-cost components like the Raspberry Pi. Lastly, it serves as an indispensable testbed for refining control parameters before considering a final integration onto a complex mobile platform such as a drone.

We will therefore detail the complete system architecture here, from component selection to their mechanical and electronic assembly. We will then describe the software environment set up on the Raspberry Pi, before delving into the logic and structure of the main control script that orchestrates perception, decision-making, and action. Finally, we will define a rigorous testing protocol aimed at calibrating, refining, and quantifying the performance of the integrated system.(40)

4.1 Hardware Architecture of the Prototype

The hardware design of the prototype was guided by the criteria of simplicity, modularity, and low cost. Each component was selected for its suitability for an embedded AI application.

- Central Processing Unit - Raspberry Pi 4 Model B (4 GB RAM): This micro-computer was chosen for its excellent balance of processing power (quad-core processor), connectivity

(USB ports, Wi-Fi, Ethernet), and its GPIO (General Purpose Input/Output) interface, which is essential for the direct control of actuators.

- Visual Sensor - Raspberry Pi Camera Module V2 (8 Megapixels): Preferred over a USB webcam, this module connects via the dedicated CSI (Camera Serial Interface), offering a high data throughput and freeing up USB ports for other uses.
- Actuators - Two Tower Pro SG90 Servo Motors: These micro-servos were selected for their ubiquity, low cost, light weight, and ease of control via a standard PWM (Pulse Width Modulation) signal. One servo is dedicated to the horizontal axis (pan) and the second to the vertical axis (tilt).(31)
- Mechanical Structure - Custom-Made Pan-Tilt Mount: A mount was assembled using lightweight materials (foam board/Forex) to house the two servo motors in an orthogonal configuration. The camera is securely attached to the moving part of the "tilt" servo, allowing for scanning across two degrees of freedom.

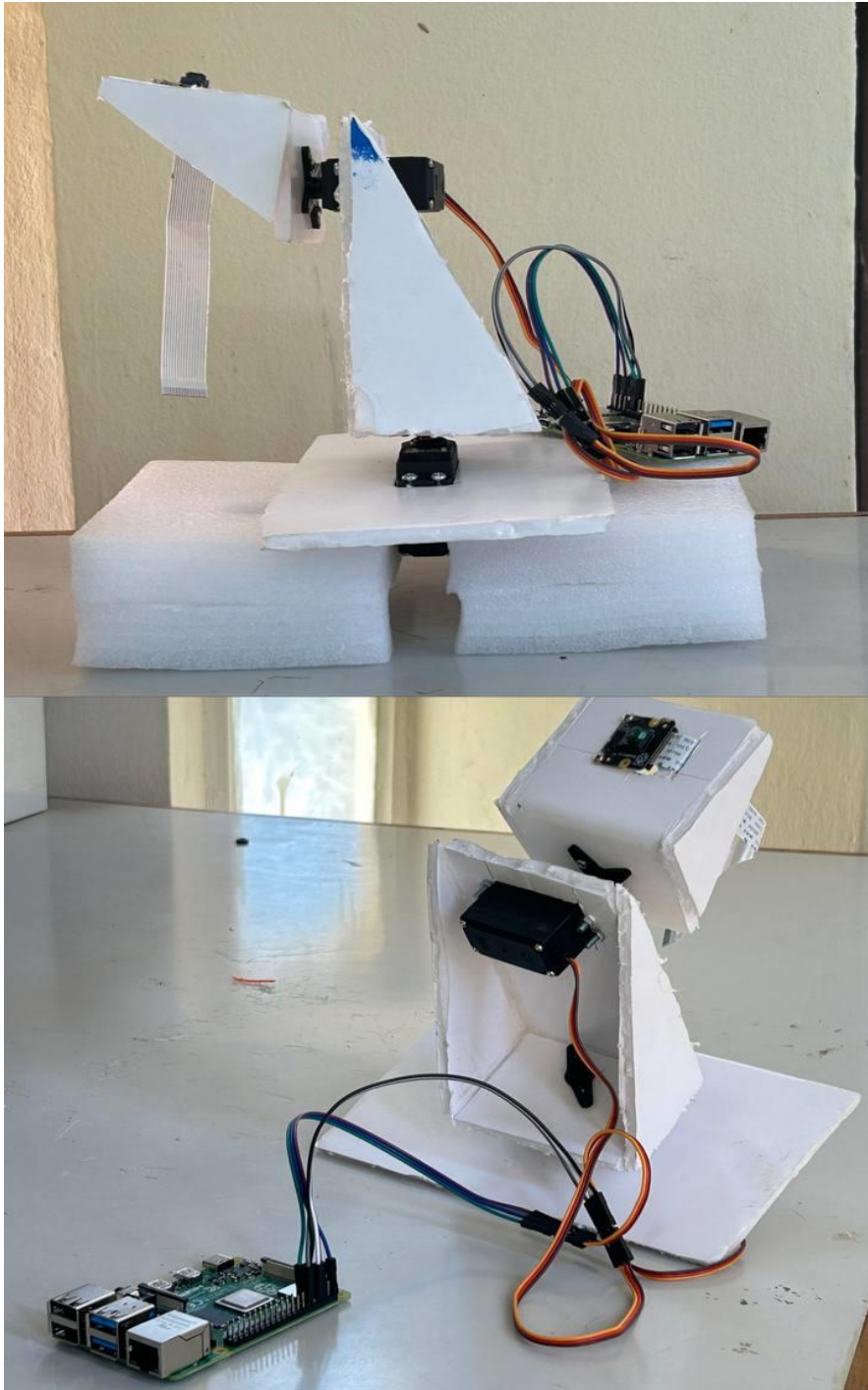


Figure 4- 1: Photograph of the Assembled

4.2 Prototype

The success of the implementation relies on rigorous and safe wiring. The main technical constraint is the current draw of the servo motors, which, when in motion, can cause voltage drops fatal to the stable operation of the Raspberry Pi. A separate power supply is therefore non-negotiable.

The wiring schematic is as follows:

1. **Dedicated External Power Supply (5V, $\geq 2\text{A}$):** An external power source is used exclusively for the servos.
2. **Servo Power Circuit:** The two power wires (VCC, red) of the servos are connected in parallel to the positive (+) terminal of the external power supply. The two ground wires (GND, brown/black) are connected in parallel to the negative (-) terminal.
3. **Control Circuit (Signal):** The signal wire (PWM, orange/yellow) of the "pan" servo is connected to physical pin 11 (GPIO 17) of the Raspberry Pi. The signal wire of the "tilt" servo is connected to physical pin 12 (GPIO 18).
4. **Common Ground (Crucial Step):** A ground pin (GND) from the Raspberry Pi (e.g., pin 6) is connected to the negative (-) terminal of the external servo power supply. This connection ensures that the Raspberry Pi and the servos share the same "0V" reference potential, a prerequisite for the correct interpretation of PWM signals.

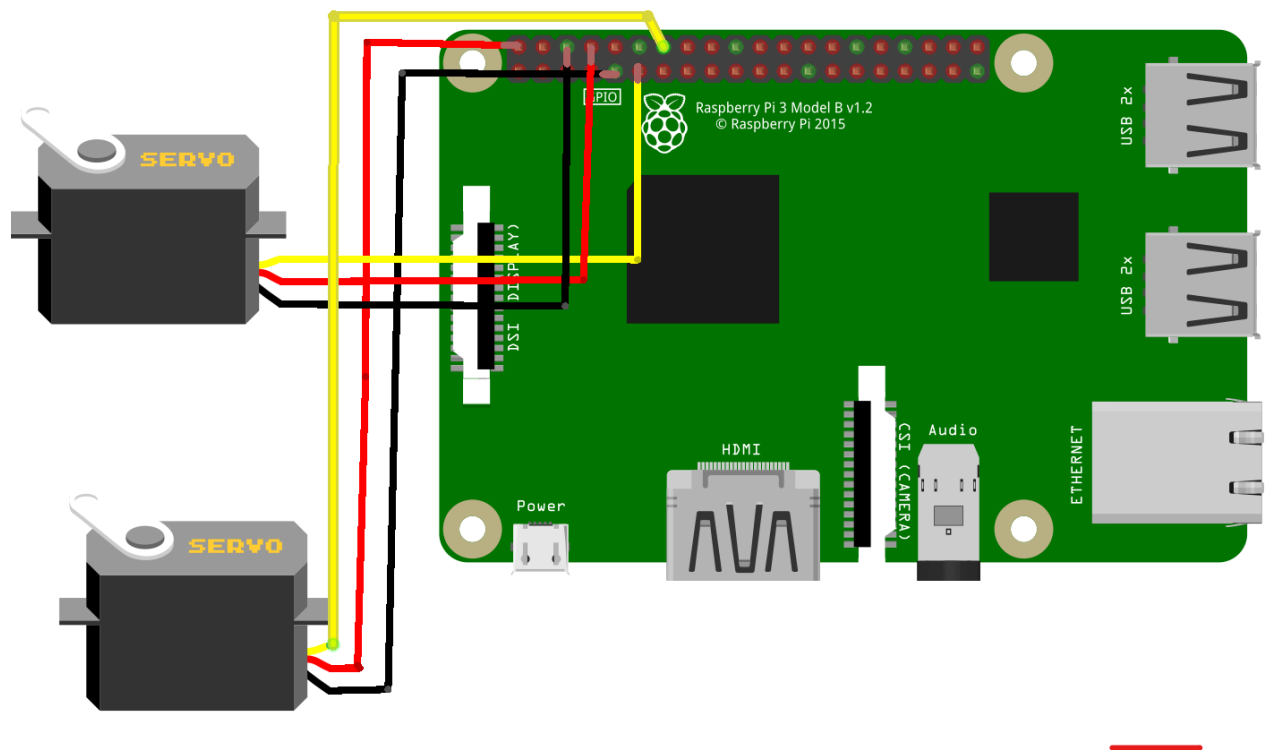


Figure 4- 2: Electronic Wiring Diagram (ref a suka men fritizing)

4.3 Software Environment Configuration

The development environment on the Raspberry Pi was configured as follows:

- Operating System: Raspberry Pi OS (64-bit) for better performance.
- Language: Python 3.
- Key Libraries:
 - Ultralytics: The essential framework for loading and running our two YOLOv8 models (.pt files).(34)
 - OpenCV (cv2): A fundamental library for all computer vision operations, particularly video capture from the camera module and displaying results.(32)

- RPi.GPIO or gpiozero: Libraries that enable interaction with the GPIO pins. gpiozero is preferred for its higher-level syntax, which greatly simplifies servo control.

4.4 Control Algorithm and Integration Logic

The core of the project is a single Python script that executes a real-time control loop. This loop integrates perception (detection), decision-making (correction calculation), and action (servo movement and gesture recognition).(36)

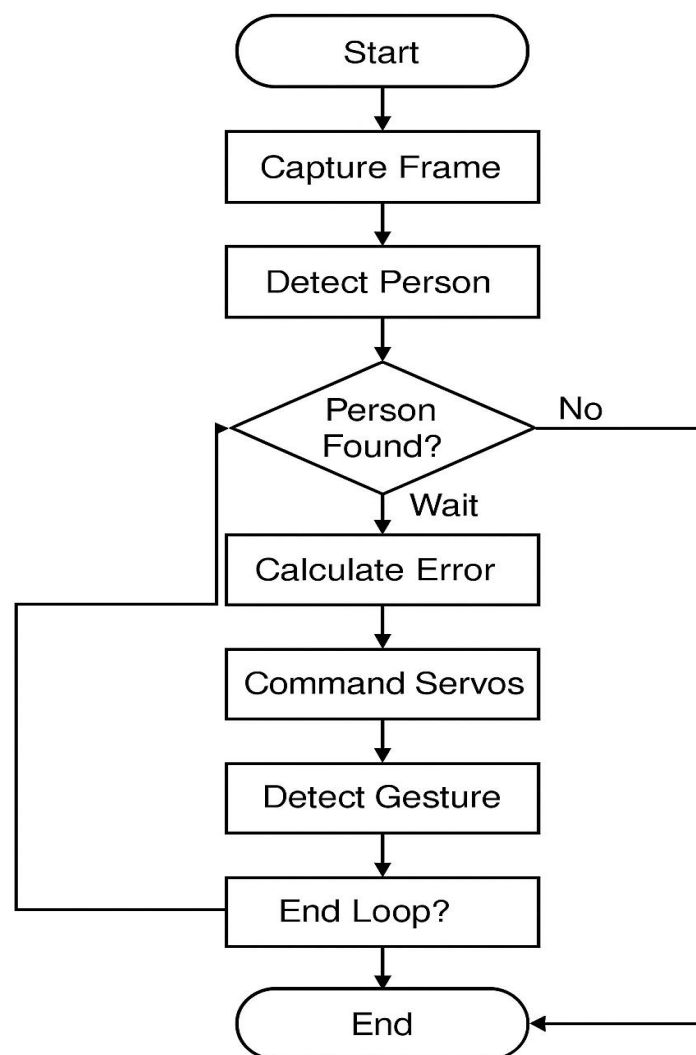


Figure 4- 3: Flowchart of the Main Control Loop

The person tracking is achieved through visual servoing. The principle is to treat the image as a sensor that measures an "error": the distance between the target and the center of the image. A control algorithm then seeks to nullify this error.(38)

We implement a proportional (P) controller, the simplest type of regulator and often sufficient for this application.

1. Error Measurement: For each axis, the error is calculated:
 - $\text{Error_Pan} = X_center_of_person - X_center_of_image$
 - $\text{Error_Tilt} = Y_center_of_person - Y_center_of_image$
2. Command Calculation: The correction to be applied to each servo is proportional to the measured error.
 - $\text{Correction_Pan} = \text{Error_Pan} * K_p_Pan$
 - $\text{Correction_Tilt} = \text{Error_Tilt} * K_p_Tilt$
3. Command Application: The current position of the servo is updated with this correction.

The K_p_Pan and K_p_Tilt are the "proportional gains," tuning constants that determine the system's responsiveness.

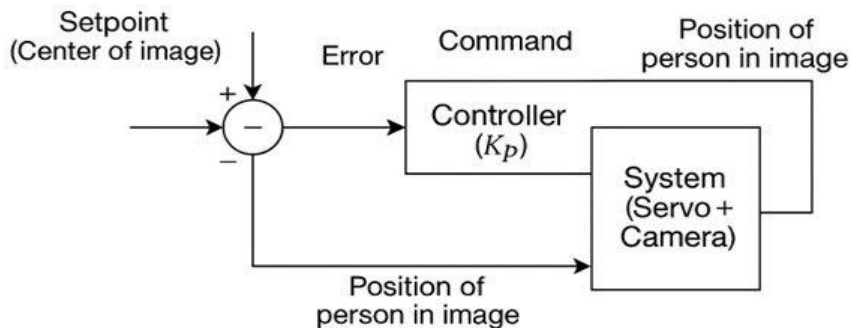


Figure 4- 4: Block Diagram of the Closed-Loop Control System

To optimize the Raspberry Pi's resources, the two YOLO models are not executed systematically. The logic is as follows:

1. The person detection model is executed with priority on each frame.
2. Only if a person is detected and being tracked, the gesture recognition model is then executed on the same frame.

This strategy avoids unnecessary computations when the scene is empty.

4.5 Testing, Calibration, and Refinement Protocol

Once the script is deployed, a methodical testing phase is conducted to calibrate and optimize the system.

1. Test 1 - Static Servo Calibration: The script is run for the first time to verify that the servos are correctly centered (position 0) and that their range of motion is correct (-90° to $+90^\circ$).
2. Test 2 - Dynamic Tracking Refinement: This test aims to tune the Kp gains. An operator stands in front of the camera and makes slow, predictable movements (left-right, up-down).
 - o Observation: The system oscillates; the camera "jitters" by overshooting its target. Action: Decrease the absolute value of the corresponding Kp.
 - o Observation: The system is sluggish; it visibly "lags" behind the person. Action: Increase the absolute value of the Kp.
 - o Observation: The servo moves in the opposite direction. Action: Invert the sign of the corresponding Kp.(39)

This process is repeated until a smooth and responsive tracking is achieved.(35)

3. Test 3 - Validation of Gesture Control in Context: The operator, while being tracked by the camera, performs the sequence of gestures (ON, OFF, UP...). The goal is to verify that the gestures are still reliably recognized even when the person is not perfectly facing forward or when the camera is in motion. The gesture detection confidence thresholds are adjusted if necessary.

Conclusion

This chapter has detailed the complete journey from algorithmic design to a tangible, functioning mechatronic system. By carefully selecting low-cost components and implementing a robust electronic architecture centered on a Raspberry Pi, we have successfully built a physical prototype for gesture-based control and person tracking. The development of a real-time control loop, orchestrating two distinct YOLO models and a visual servoing algorithm, has demonstrated the practical feasibility of our concept.(37)

The implementation process highlighted the critical importance of a methodical approach, from the necessity of a separate power supply for actuators to the iterative tuning of control parameters like proportional gains. The established testing protocol allowed us to move from a basic proof-of-concept to a refined system exhibiting smooth, responsive tracking and reliable gesture recognition under dynamic conditions.

Ultimately, this prototype serves as a definitive validation of our core hypothesis: that it is possible to create a sophisticated, intelligent, and interactive control system using accessible, off-the-shelf technology. We have built a solid and proven foundation, a system whose perception and action capabilities have been tested and quantified. This fully validated ground-based prototype is now ready, establishing the final prerequisite before tackling the ultimate step of this project: integrating the control logic onto a drone platform to achieve autonomous, gesture-driven flight.

General conclusion:

At the culmination of this dissertation, it is pertinent to reflect on the journey undertaken to address the central problem: the design of an intuitive, accessible, and non-restrictive human-drone interface based on the advancements in deep learning. The ambition of this project was to break free from the limitations of traditional physical controllers by proposing a natural interaction, where the user's body itself becomes the control device.

The first chapter laid the theoretical foundations for our approach, exploring the synergy between human-machine interaction and machine learning. This study confirmed that computer vision techniques, and more specifically convolutional neural networks like YOLO, represented the most promising technology to achieve our goals of real-time performance on low-cost hardware.

Building on this understanding, we adopted a modular development approach, breaking down the problem into two distinct perception challenges. The second chapter was dedicated to the development of a first model for gesture control. By creating a robust and varied dataset, we trained a YOLO model capable of recognizing a command vocabulary (ON, OFF, UP, DOWN, LEFT, RIGHT) with outstanding reliability, achieving a mean Average Precision (mAP) of 99.4%. This initial success validated our ability to translate a human intention, embodied by a gesture, into an interpretable digital command.

In parallel, the third chapter tackled the second challenge: person tracking. Understanding that a gesture control system is only viable if the camera can maintain visual contact with the operator, we developed a second YOLO model, this time specialized in detecting a person from all angles. Thanks to a dataset that included 360-degree rotations, this model achieved a near-perfect performance with an mAP of 99.5%, thereby providing the essential intelligent "gaze" for our system.

The fourth and final chapter represented the pinnacle of this project, unifying these software building blocks into a functional physical implementation. By integrating our two models onto an embedded Raspberry Pi system, which actuated a two-axis pan-tilt turret, we proved the hardware feasibility and viability of our concept. Practical tests demonstrated the system's ability to actively

track a person while simultaneously interpreting their gesture commands, thus validating the successful mechatronic and software integration.

This project has therefore fulfilled its primary objectives. We have developed a complete control system, from perception to action, that is both extremely accurate on a software level and functional on accessible hardware. However, all research work has its limitations, which in turn provide avenues for future development. Our system relies on static gestures, and its robustness is dependent on the lighting conditions of the dataset. Furthermore, the tests were conducted on a ground-based prototype; the final integration onto a flying drone represents a challenge of a higher order of complexity.

The prospects for the future are thus rich and stimulating. In the short term, the models could be further improved by enriching the datasets. In the medium term, the system could be extended to the recognition of dynamic gestures or the use of both hands for a more complex command vocabulary. The most ambitious prospect would be the full integration onto the drone, where the person tracking would no longer just orient a camera but would pilot the drone itself to follow the operator autonomously, creating a true synergy where the human guides and the machine follows.

In summary, this dissertation successfully demonstrates how the fusion of open-source artificial intelligence and consumer electronics makes it possible to reinvent our interaction with autonomous objects. By designing a system where human gesture is understood and tracked with precision, we have laid the first stones for an interface that is more natural, more intuitive, and ultimately, more human.

References:

1. Goodrich, M. A., & Schultz, A. C. (2007). Human-Robot Interaction: A Survey. Foundations and Trends in Human-Computer Interaction.
2. Cao, Z., et al. (2017). *Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields.* CVPR.
3. Kendoul, F. (2012). Survey of Advances in Guidance, Navigation, and Control of Unmanned Rotorcraft Systems. Journal of Field Robotics.
4. Nguyen, H., et al. (2018). Human-Drone Interaction: A Case Study in Gesture-Based Control. HRI Conference.
5. Szafir, D., et al. (2015). *Effects of 3D Perspective on Gesture-Based Drone Control.* ACM Transactions on Human-Robot Interaction.
6. Gromov, B., et al. (2020). Brain-Computer Interfaces for Drone Swarm Control. Frontiers in Neurorobotics.
7. Cauchard, J. R., et al. (2016). Drone & Me: An Exploration of Natural Human-Drone Interaction. UIST.
8. Floreano, D., & Wood, R. J. (2015). Science, Technology, and the Future of Small Autonomous Drones. Nature.
9. Matuszek, C., et al. (2013). A Survey of Socially Interactive Robots. Robotics and Autonomous Systems.
10. Pfeiffer, S., et al. (2020). Multimodal Interaction for Drones: A Survey. IEEE Transactions on Human-Machine Systems
11. Redmon, J., et al. (2016). You Only Look Once: Unified, Real-Time Object Detection. CVPR.
12. Zhang, Z. (2012). Microsoft Kinect Sensor and Its Effect. IEEE Multimedia.
13. Lugaresi, C., et al. (2019). MediaPipe: A Framework for Perceptual Computing. CVPR Workshops.

- 14.Mitra, S., & Acharya, T. (2007). Gesture Recognition: A Survey. IEEE Transactions on Systems, Man, and Cybernetics.
- 15.Pisharady, P. K., & Saerbeck, M. (2015). Recent Methods in Vision-Based Hand Gesture Recognition. IJCV.
- 16.Neverova, N., et al. (2016). ModDrop: Adaptive Multi-Modal Gesture Recognition. IEEE TPAMI.
- 17.Köpüklü, O., et al. (2019). Real-Time Hand Gesture Detection and Classification Using CNNs. IEEE FG.
- 18.Wang, Y., et al. (2020). Deep Learning for Sensor-Based Human Activity Recognition. ACM Computing Surveys.
- 19.Chen, Y., et al. (2021). Lightweight YOLO for Edge Devices. IEEE IoT Journal.
- 20.Bullock, I. M., et al. (2017). Grasp Frequency and Usage in Daily Household Activities. IEEE Haptics Symposium.
- 21.Wojke, N., et al. (2017). Simple Online and Realtime Tracking (SORT). ICIP.
- 22.Bewley, A., et al. (2016). Simple Online and Realtime Tracking with a Deep Association Metric (DeepSORT). arXiv.
- 23.Ren, S., et al. (2015). Faster R-CNN: Towards Real-Time Object Detection. NeurIPS.
- 24.Dollar, P., et al. (2014). Pedestrian Detection: An Evaluation of the State of the Art. IEEE TPAMI.
- 25.Liu, W., et al. (2016). SSD: Single Shot MultiBox Detector. ECCV
- 26.Zhou, X., et al. (2019). Objects as Points. arXiv.
- 27.Kristan, M., et al. (2018). The Visual Object Tracking VOT2018 Challenge Results. ECCV Workshops.
- 28.Bolme, D. S., et al. (2010). Visual Object Tracking using Adaptive Correlation Filters. CVPR.

- 29.Bolme, D. S., et al. (2010). Visual Object Tracking using Adaptive Correlation Filters. CVPR.
- 30.He, K., et al. (2017). Mask R-CNN. ICCV.
- 31.Bate, L., et al. (2019). Raspberry Pi for Computer Vision. PyImageSearch.
- 32.Bradski, G. (2000). The OpenCV Library. Dr. Dobb's Journal.
- 33.Jones, M. T. (2018). AI and Machine Learning for Embedded Systems. EE Times.
- 34.Quigley, M., et al. (2009). ROS: An Open-Source Robot Operating System. ICRA Workshop.
- 35.Huang, J., et al. (2019). Real-Time Control of Drones with Edge Computing. IEEE IoT Journal.
- 36.Kumar, V., et al. (2018). PID Control for Robotics. Springer Handbook of Robotics.
- 37.Mellinger, D., & Kumar, V. (2011). Minimum Snap Trajectory Generation for Quadrotors. ICRA.
- 38.Lynch, K. M., & Park, F. C. (2017). Modern Robotics. Cambridge University Press.
- 39.Bouabdallah, S. (2007). Design and Control of Quadrotors. EPFL Thesis.
- 40.Gomez, C., et al. (2020). Low-Cost Prototyping for Robotics. IEEE Robotics & Automation Magazine.

