



University of Reading  
Department of Computer Science

# Computer Vision for Multi-Object tracking

Uddeeph Dasari

*Supervisor:* Dr. James Ferryman

A report submitted in partial fulfilment of the requirements of  
the University of Reading for the degree of  
Master of Science in *Data Science and Advanced Computing*

10<sup>th</sup> September 2024

## Declaration

I, Uddeepth Dasari, of the Department of Computer Science, University of Reading, confirm that this is my work and that figures, tables, equations, code snippets, artworks, and illustrations in this report are original and have not been taken from any other person's work, except where the works of others have been explicitly acknowledged, quoted, and referenced. I understand that failing to do so will be considered a case of plagiarism. Plagiarism is a form of academic misconduct and will be penalised accordingly.

I give consent to a copy of my report being shared with future students as an exemplar.

I give consent for my work to be made available more widely to members of UoR and the public with an interest in teaching, learning and research.

Uddeepth Dasari  
10 September 2024

## Abstract

The project aims to develop an advanced system for real-time multi-object tracking using DeepSORT and YOLOv5. Its purpose is to enhance accuracy in identifying and tracking multiple objects in video streams. YOLOv5 detects objects. DeepSORT tracks them. This enables real-time monitoring in autonomous vehicles, surveillance, and robotics. This fusion powers critical applications, enhancing safety and efficiency across industries. Its precision and speed are vital for modern tech. It keeps object identities across multiple frames. Also, the project involves a complex visualisation module. It must generate heat maps for movement analysis. It must also provide metrics to evaluate the system's performance under different conditions. The results show a big boost in tracking accuracy. This is true for high-density, fast-moving object scenarios. It proves the integrated system works well. These advancements lay a strong base for future improvements in real-time object tracking. They could benefit many industries. The project's results suggest good paths for future research. They are to optimise tracking algorithms and improve system robustness.

Code Link: [https://gitlab.act.reading.ac.uk/jm837123/csmpr21\\_31837123\\_code.git](https://gitlab.act.reading.ac.uk/jm837123/csmpr21_31837123_code.git)

Keywords: Multi-Object Tracking, YOLOv5, DeepSORT, Real-Time Object Detection, Thermal and Visible Data, Autonomous Vehicles, Surveillance, Robotics, Tracking Accuracy

## Acknowledgements

I would like to extend my heartfelt thanks to my supervisor, Dr. James Ferryman, for his invaluable guidance and support throughout this project. Their expertise and feedback were crucial to the success of this research. I also appreciate the Department of Computer Science for providing the resources and facilities essential for completing this work.

# Table of Contents

| Section Title   | Page Number |
|---|-------------|
| <b>1. Introduction</b>                                  |             |
| 1.1 Background  | 1           |
| 1.2 Problem Statement                                   | 2           |
| 1.3 Aims and Objectives                                 | 2           |
| 1.4 Significance and Relevance                          | 3           |
| 1.5 Solution Approach                                   | 3           |
| 1.6 Summary of Contributions and Achievements           | 4           |
| 1.7 Organization of the Report                          | 5           |
| <b>2. Literature Review</b>                             |             |
| 2.1 Overview of Multi-Object Tracking                   | 6           |
| 2.2 YOLO and Real-Time Object Detection                 | 6           |
| 2.3 DeepSORT and Object Tracking                        | 7           |
| 2.4 Challenges in Multi-Object Tracking                 | 7           |
| 2.5 Thermal and Visual Data Fusion for Tracking         | 8           |
| 2.6 Summary of Findings and Gaps                        | 9           |
| 2.7 Visual Diagrams and Approaches                      | 9           |
| <b>3. Methodology</b>                                   |             |
| 3.1 Data Collection and Preprocessing                   | 14          |
| 3.2 Training YOLO on Thermal and Visible Data           | 17          |
| 3.3 Integration of YOLO with DeepSORT                   | 18          |
| 3.4 System Design and Architecture                      | 18          |
| 3.5 Metrics Definition and Equations                    | 21          |
| 3.6 Performance Metrics and Evaluation                  | 22          |
| 3.7 Tools and Technologies                              | 23          |
| <b>4. Results</b>                                       |             |
| 4.1 Quantitative Results: MOTA, MODA, and CLEAR Metrics | 24          |
| 4.2 Comparison with Existing Systems                    | 28          |
| 4.3 Qualitative Results: Visualizations and Heatmaps    | 30          |
| 4.4 Sequence Results and Dataset Discussion             | 33          |
| 4.5 Summary of Results                                  | 34          |
| <b>5. Discussion and Analysis</b>                       |             |
| 5.1 Analysis of Tracking Performance                    | 35          |
| 5.2 Evaluation of Data Fusion Techniques                | 35          |
| 5.3 Strengths and Limitations of the System             | 37          |
| 5.4 Implications for Real-Time Applications             | 38          |
| 5.5 Illustration of Discussion with Visuals             | 39          |
| <b>6. Conclusions and Future Work</b>                   |             |
| 6.1 Conclusions   | 40          |
| 6.2 Future Work   | 40          |
| <b>7. Reflections</b>                                   | 42          |

# List of Figures

| Fig No | Name of the figure                 | Page No |
|--------|------------------------------------|---------|
| 2.1    | Thermal and Visible Data Fusion    | 12      |
| 2.2    | Real-Time Processing Flow          | 13      |
| 3.1    | Visible spectrum image             | 16      |
| 3.2    | Thermal image                      | 17      |
| 3.3    | System Architecture Diagram        | 21      |
| 4.1    | Visualization of CLEAR metrics     | 31      |
| 4.2    | Heatmap Showing Occlusion Handling | 31      |
| 4.3    | Comparative Analysis of Models     | 32      |
| 4.4    | Comparison with EUMARS             | 32      |
| 4.5    | Simulated Scenario Analysis        | 32      |

# List of Tables

| Table No | Table Name                               | Page No |
|----------|--|---------|
| 4.1      | Training with YOLOv5 using CLEAR Metrics | 25      |
| 4.2      | Integration of YOLOV5 with DeepSORT      | 27      |
| 4.3      | Comparative Analysis                     | 28      |
| 4.4      | Combined Additional                      | 29      |
| 4.5      | Comparison with EUMARS                   | 30      |

# List of Abbreviations

- YOLO - You Only Look Once
- MOTA – Multi-Object Tracking Accuracy
- MODA – Multi-Object Tracking Accuracy
- MOT - Multi-Object Tracking
- MOTP – Multi-Object Tracking Precision
- CLEAR – Classification and Evaluation of Events, Activities, and Relationships
- LiDAR – Light Detection and Ranging
- RADAR - Radio Detection and Ranging
- FP - False Positives
- FN – False Negatives
- IDS – Identity Switches
- GT – Ground Truth Objects
- DeepSORT - Deep Simple Online and Real-time Tracking
- EUMARS - European Multi-Authority Border Security



# Chapter 1

## Introduction

### 1.1 Background

Real-time tracking of multiple objects revolutionizes computer vision. This crucial technology empowers autonomous vehicles, robots, and security systems. MOT tracks multiple targets at once. This is crucial in many fields. It highlights the need for fast, accurate object detection and monitoring. For example, self-driving cars must always watch for other cars, people, and obstacles. This is crucial for safe driving. Also, surveillance systems must track multiple people or objects accurately. This is vital for public safety. In robotics, MOT lets machines interact with their surroundings. It helps them respond to real-time changes.

Deep learning has improved multi-object tracking systems. YOLO revolutionizes object detection. This groundbreaking model scans images just once, swiftly identifying multiple objects. Its speed and accuracy make it a game-changer in computer vision. YOLO processes images incredibly fast while keeping high accuracy. Unlike older object detection systems, YOLO uses a single pass to process an image. It can detect multiple objects in that pass. This is crucial for applications needing quick responses. YOLO can detect multiple objects in one frame. Thus, it's valuable for industries that require immediate object detection.

YOLO is great for detecting objects. However, it struggles with tracking them in complex settings. Multi-object tracking systems struggle with overlapping, fast-moving, or partially hidden objects. They also find it tough in low-visibility conditions, like at night or in fog. Most of these systems use only visible spectrum cameras. These cameras fail in poor lighting.

This project combines thermal imaging with visible data. It will create a better tracking system to address these challenges. The system merges thermal and visible data to track objects in low visibility. It uses YOLOv5 for detection and DeepSORT for tracking across frames. Adding thermal data boosts performance when visible-light cameras falter.

## 1.2 Problem Statement

Despite recent advancements in multi-object tracking, several critical challenges remain. Multi-object tracking faces key challenges. One major issue is tracking objects in crowded or dynamic settings. When objects are close or move unpredictably, systems can lose track of them. This is worse when objects overlap or hide behind each other, causing occlusion. Such situations can lead to wrong identifications or missed detections.

Another challenge is the heavy reliance on visible spectrum data. Most tracking systems depend on visible light, making them fail in poor visibility. Low-light environments, like at night or in fog, drastically reduce accuracy. Here, cameras might not see objects, leading to tracking errors or mix-ups. This is critical for security and autonomous vehicles, where accurate tracking is vital.

This project aims to overcome these issues. It combines YOLOv5 with DeepSORT for better detection and tracking. Moreover, it adds thermal data fusion to improve tracking in low-visibility conditions. By using both thermal and visible light data, the system can detect objects by their heat. This method works even in darkness or fog. The goal is to improve tracking in crowded, low-visibility places. It should be more reliable than traditional systems.

## 1.3 Aims and Objectives

This project aims to build a real-time multi-object tracking system. It will combine YOLOv5's detection with DeepSORT's tracking. The system will track multiple objects in video frames. It will work well even in low visibility or crowded scenes. Adding thermal and visible data fusion will boost its performance. This enhancement allows tracking in conditions where standard systems fail.

The specific objectives of this project are:

- Create a dataset with thermal and visible images for better tracking.
- Train the YOLOv5 model on this dataset for real-time object detection, even in bad conditions.
- Combine YOLOv5 with DeepSORT to track objects across frames and keep their identities. Develop a method to merge detection and tracking data.
- Use MOTA and MODA to evaluate the system. MOTA checks tracking accuracy, while MODA looks at detection precision.
- Compare the new system with others. Show it is faster, more accurate, and more reliable in tough conditions.

## 1.4 Significance and Relevance

This project is important. It focuses on tracking multiple objects in real time. From public safety to robotics, this innovation tracks vital data. Self-driving cars rely on it to spot vehicles, people, and hazards. Passengers and pedestrians alike depend on its precision. Beyond transportation, its reach extends to various fields, enhancing safety and efficiency. As the technology evolves, so do its applications, shaping our future. Similarly, public safety systems must monitor multiple subjects to prevent threats.

The project stands out due to its unique method. It combines thermal and visible data. This approach fixes a major flaw in current systems that rely only on visible data. These systems struggle in low visibility, like at night. Our system, however, continues to track objects in darkness or bad weather. This is thanks to the added thermal data.

Speed is another key feature of our system. Applications like autonomous vehicles and public surveillance need quick processing. A slow system can't react in time. For example, a car must quickly respond to a pedestrian. Similarly, a surveillance system must track people in real-time to catch potential threats. Our project aims to balance speed and accuracy, making it ideal for real-time needs.

## 1.5 Solution Approach

This section describes the strategies formulated during the project to overcome the obstacles of real-time multi-object tracking. The solution employs both the YOLOv5 model for real-time object detection and DeepSORT for tracking objects from frame to frame.

An organized system can detect and track objects in real-time and under more challenging conditions such as those where visibility is limited and could otherwise interfere with the traditional systems. When implementing the combination of YOLOv5 and DeepSORT its purpose is to make the processes of detecting objects in frames and tracking them over subsequent frames in addition to preserving their identities both fast and accurate as well as robust to such factors as occlusion or when some of the object leaves the field of vision temporarily.

The development environment for this project entails the use of Python as the main programming language with the deep learning framework being PyTorch. PyTorch is selected due to its compatibility with different architectures, and pre-trained models such as YOLOv5, which well fits the implementation of a real-time object detection use case. Image and video processing is done by the OpenCV library and to track object identities across frames DeepSORT algorithm is employed.

In this solution, YOLOv5 oversees object detection per frame, while DeepSORT uses object detection results of multiple frames to track objects. The primary novelty in this project is the combination of thermal and visible data because the system would be able to track objects in low light or deteriorating weather conditions where conventional visible spectrum cameras would not be useful.

This approach provides enough assurances that the detection and tracking can be reliable making it suitable for use in areas like self-driving cars, surveillance systems and robotics where speed and accuracy are essential.

## 1.6 Summary of Contributions and Achievements

The following are some of the outstanding benefits of this project in the field of multi-object tracking. Among these, they include real-time object detection using YOLOv5 which was a success. On the other hand, DeepSORT has an excellent tracking algorithm by applying a tracking algorithm side by side with YOLOv5 which provides the simplicity of detecting multiple objects within a short period. Due to their effective combination, this can be especially beneficial in environments in which the objects are likely to interpenetrate or may often translate in unpredictable manners. It is even noteworthy that, in such poor conditions, the tracking accuracy of the system remains extremely high.

The other accomplishment is the integration of thermal and visible data sets, an advancement that provides an enhancement of digital maps. This goes a long way in improving the tracking of objects in scenarios in which other tracking systems would not manage to get it done. Given that object detection is based on thermal data, the system can track objects even under ill-illumination such as at night or in cases where the area is foggy. This makes the system more flexible as well as allows the system to work in various conditions as compared to other tracking systems which operate in the visible spectrum of light.

Finally, for real-time operation, the system has been adjusted in such a way that it can run video streams suitable for practical applications such as autonomous driving and security camera systems. This is already apparent if one investigates the system's tracking results based on the standard markers including MOTA and MODA that evidence the efficiency of the system in object identification aspect as well as while maintaining object identity numbers across frames.

## 1.7 Organization of the Report

The organization of this report is such that it will cover all aspects of the project right from the background and statement of the problem to the results and conclusion section. The report is organized as follows:

- **Introduction:** This chapter briefly introduces the project and discusses major issues of MOT and how this project can solve them. It also presents the findings and contributions of the project.
- **Literature Review:** This chapter is a literature review on literature regarding multi-object tracking with specific emphasis on the YOLO detection model, DeepSORT tracking algorithm, and a multimodal approach that incorporates thermal and visible data. The literature review section points out the research gaps that prevail in the current carry-out and shows how this project aims to close those gaps.
- **Methodology:** In this section, the process of pre-processing of the data set, training of the YOLOv5 model and its incorporation to DeepSORT through the explanation given. It also describes system design, data integration techniques, and the way, in which the system's effectiveness is evaluated.
- **Results:** Based on these criteria, this chapter describes the results of an evaluation of the system, thus providing the values of MOTA and MODA as well as examples of visual tracking. It gives a plan of the system's ability to perform and get accurate results in different scenarios.
- **Discussion:** This chapter also compares the performance of the system with other existing multi-object tracking systems. In addition, it provides analysis of areas of strength, areas of weakness, and areas which have potential for more advancement.
- **Conclusion and Future Work:** This section provides a general conclusion based on the result obtained and highlights the contributions of the project and proposed areas for further research and development.
- **Reflections:** This chapter focuses on the difficulties experienced during the project as well as some of the lessons gained. It also argues about the competencies acquired during development.

# Chapter 2

## Literature Review

### 2.1 Overview of Multi-Object Tracking

MOT is to find and follow multiple objects in a video across frames. Its importance spans many fields, like autonomous vehicles and robotics. MOT's key goal is to detect and track objects as they move through a scene. Despite overlaps, direction changes, or size variations, it must maintain its identity.

Historically, MOT relied on simpler methods such as Kalman filtering and optical flow. However, deep learning models have improved MOT systems. They are now faster and more accurate. This enables practical real-time applications.

### 2.2 YOLO and Real-Time Object Detection

YOLO has significantly advanced real-time object detection due to its speed and efficiency. Unlike older models, which processed images in separate regions, YOLO processes the entire image in one pass. This significantly improves detection speed. However, while YOLO is faster than other models like Faster R-CNN, it comes with certain trade-offs. For instance, YOLO struggles to detect smaller objects in complex, crowded scenes. Its grid-based system may cause smaller objects to be missed if they don't dominate the grid cell. Also, overlapping or too-close objects can cause misclassification.

YOLO can detect objects in real time. However, it does not retain information across frames. So, it is not suitable for multi-object tracking by itself. The system struggles to maintain object identity over time. This is critical in MOT scenarios.

In this project, the YOLOv5 model is integrated with DeepSORT to improve tracking across frames. YOLO detects objects. DeepSORT extends it by keeping track of object identities over time, even when they overlap or briefly disappear. Also, the project aims to improve YOLO's detection in low-visibility conditions, like at night or in fog. It will do this by using thermal-visual data fusion. Visible-spectrum cameras alone may fail in these conditions. This combination helps fix YOLO's flaws. It tracks smaller and occluded objects more reliably in various environments.

## 2.3 DeepSORT and Object Tracking

DeepSORT builds on the SORT algorithm. SORT uses motion predictions to track objects. DeepSORT adds a deep learning component. It uses appearance-based features to maintain object identities more accurately. This helps, especially when objects temporarily disappear or are occluded by other objects. This approach improves tracking in crowded places. There, objects often occlude and need re-identification.

However, DeepSORT has its limitations. The system can struggle with objects that change a lot due to lighting, orientation, or camera angles. For example, an object may be misidentified after it changes appearance when it moves in or out of shadows, or when the view angle changes. Also, DeepSORT may struggle to track objects in crowded areas with frequent, close interactions.

This project uses thermal-visual data fusion to improve tracking accuracy. It aims to address these challenges. Combining thermal and visible-spectrum data makes the system more robust. It helps in cases of changing light or occluded objects. Thermal data shows the heat signatures of objects. It can help identify them even if their appearance changes in the visible spectrum. This approach lets the system handle more complex environments. These include areas with dense crowds and frequent object interactions.

## 2.4 Challenges in Multi-Object Tracking

Despite recent advancements, MOT continues to face several persistent challenges. A major challenge is occlusion. Objects may overlap or hide behind others. This makes it hard to track their identities as they move through the scene. This issue is common in crowded places, like busy streets or large gatherings. There, objects constantly obstruct each other.

Another key challenge is scale variation. Objects look smaller or larger depending on their distance from the camera. This variation can confuse tracking systems, leading to incorrect predictions of object movements. Many MOT systems struggle with environments where objects often move in and out of the frame or change direction quickly.

Moreover, MOT systems are often optimized for ideal conditions, such as good lighting and clear weather. Tracking systems must work in poor lighting, like at night or indoors, and in severe weather, like fog and rain. Traditional visible-spectrum cameras struggle in these situations, making accurate object tracking more difficult.

This project addresses these challenges by utilizing thermal and visible data fusion. Thermal cameras use heat to detect and track objects. They work in low-light or poor visibility conditions. By combining this data with visible-spectrum data, the system can track objects better. It works even in poor lighting or when objects are blocked. Merging thermal and visible data helps with scale variation. Thermal data is less affected by distance or lighting changes. It provides better tracking across different environments.

## 2.5 Thermal and Visual Data Fusion for Tracking

Many researchers have turned to thermal imaging. It can overcome the limits of traditional tracking systems. These systems rely solely on visible-spectrum cameras. Thermal cameras detect heat signatures. They can track objects in low-light or poor-visibility conditions, like at night or in bad weather. Visible-spectrum cameras give detail and colour in daylight. But they struggle in poor or no light. Thermal cameras are not affected by lighting. So, they are useful for object tracking in tough conditions.

Combining thermal and visible data can improve object-tracking systems. It will make them more robust in different conditions. However, fusing these two modalities presents challenges, particularly data alignment. Thermal and visible cameras often have different resolutions or fields of view. This combines their outputs with cohesion and precision. Also, real-time data fusion can be costly. It may slow the tracking system.

This project fuses thermal and visible data. It aims to boost performance in low-visibility conditions. The project aligns and synchronises the thermal and visible data. This enables better object detection and tracking in various conditions. The system optimizes real-time data fusion. It will not slow down with the extra work. This makes the system more dependable in the real world. Lighting and weather conditions shift unpredictably in that location.



## 2.6 Summary of Findings and Gaps

The reviewed literature shows that great progress has been made in object detection and tracking. YOLO is known as one of the best models for real-time object detection. DeepSORT is a powerful tool for tracking objects across frames using motion and appearance data. Together, these models offer significant advantages for MOT.

However, several gaps remain in the research. Many MOT systems, like YOLO and DeepSORT, struggle to track objects in low visibility or when they are occluded. Some research has looked at fusing thermal and visible data. However, there is still much room for improvement. This includes real-time performance and aligning the two data types.

This project will close these gaps. It will improve the system's performance in low-visibility areas by fusing thermal and visual data. By using these methods, the system can track objects better in both day and night. YOLOv5 and DeepSORT work together to keep object identities over time. This holds even when objects are occluded or move in and out of view. This project aims for real-time performance. It allows use in apps, like surveillance and self-driving cars, where speed and accuracy are vital.

## 2.7 Visual Diagrams and Approaches

In this section, we delve into the architecture and approach adopted in the project, providing a detailed explanation of the components and flow of the multi-object tracking system. The system is designed to combine state-of-the-art object detection and tracking methods with thermal-visual data fusion to improve performance in various challenging scenarios.

### 2.7.1 Overview of the System Architecture

The object tracking module and the object detection module make up the two main parts of the MOT system's architecture. Together, these two modules provide precise object detection and reliable object tracking across frames in both visual and thermal imagery.

1. **Object Detection Module (YOLOv5):** Finding items in each frame is the initial stage in the MOT system. The YOLOv5 model, which is renowned for its real-time object identification abilities, is used by the project for this reason. After processing the incoming photos (both visual and thermal), YOLOv5 recognizes different items in each frame.

**YOLOv5 Method:** Each cell in the grid-like representation of the image forecasts bounding boxes and the likelihood of objects.

Because it uses a single-shot method of object detection—processing the entire image in a single pass rather than scanning specific regions—YOLOv5 is quick and effective.

The system can detect objects with a fair degree of accuracy and speed by utilizing YOLOv5. But as was previously mentioned, YOLO by itself is unable to provide the temporal consistency required for MOT; this is where the tracking module is useful.

2. **Object Tracking Module (DeepSORT):** After objects are identified, their motion is tracked over a series of frames using the DeepSORT algorithm. By adding appearance features retrieved using a deep learning network, DeepSORT improves upon the Simple Online and Realtime Tracking method.

**DeepSORT Procedure:** It uses the bounding box coordinates provided by YOLOv5 as input. Using a Kalman filter, DeepSORT predicts each object's future location to track it. Additionally, each object has appearance descriptors generated for it, which aids in the system's ability to identify and recognize objects even when they are momentarily obscured or encounter other objects. When combined, YOLOv5 and DeepSORT offer both temporal and spatial tracking, guaranteeing that objects are identified in every frame and that their motions are precisely tracked over time.

### 2.7.2 Thermal-Visual Data Fusion

It combines visual and thermal data to improve object tracking and detection reliability. Thermal data, which records the heat signatures of objects, is vital for recognizing things in poor vision conditions, such as at night or in fog. In contrast, during daytime hours, visible-spectrum data offers more detail and colour information.

#### Method of Data Fusion

- **Alignment:** Aligning images from various sensors is a major difficulty in thermal-visual data fusion because of the sensor's fields of view and resolutions very often. To overcome this, the project aligns the images geographically and normalizes and calibrates the images from visible and thermal cameras, ensuring that relevant objects are recognized consistently across the two modalities.
- **Fusion Technique:** Before sending the data to the YOLOv5 model for object detection, the system first fuses the features from both visual and thermal images in a middle stage. This makes it more likely that things found in the visual data will also be confirmed by the thermal data, and vice versa.

The accuracy of object detection and tracking in low-visibility situations, including nighttime surveillance or foggy circumstances, is significantly increased by this fusion strategy. Combining the benefits of visual and thermal data makes the system more resistant to changes in illumination and surroundings.

### 2.7.3 Real-Time Processing Considerations

The project intends to accomplish real-time processing to make it suitable for practical applications, such as surveillance, autonomous navigation, and border security. Consequently, several optimizations have been implemented to lower processing overhead:

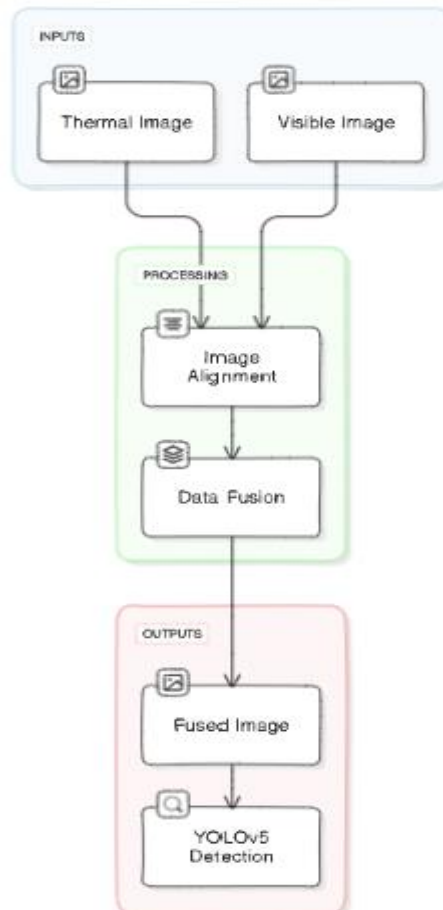
- **Effective Model Architecture:** The lightweight architecture of YOLOv5 enables quick processing without appreciably sacrificing accuracy.
- **Parallel Processing:** To maintain system responsiveness even when processing two modalities of data, thermal and visual data streams are handled independently before fusion.
- **Optimized Fusion Algorithms:** By streamlining the data fusion process to reduce latency, it is possible to prevent the system from being slowed down by the extra effort of merging visual and thermal data.

### 2.7.4 Diagram Explanation

The following diagrams show the procedures involved in the object identification, tracking, and data fusion processes to visualize the system architecture and flow. Each diagram describes a particular system component and shows how visual and thermal data are combined for improved object-tracking performance:

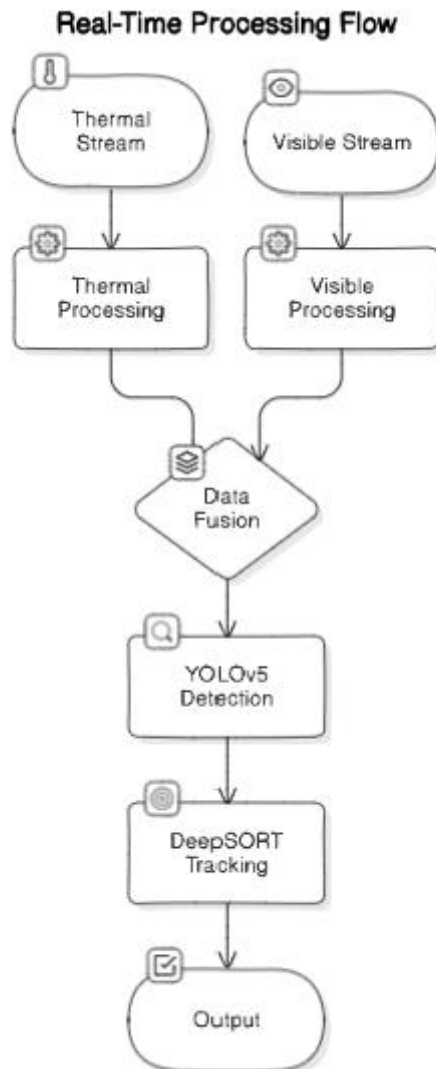
1. **Thermal and Visible Data Fusion Diagram:** This diagram demonstrates how thermal and visible images are aligned and fused before being passed through the detection and tracking pipeline.

**Thermal and Visible Data Fusion Diagram**



**Fig 2.1 Thermal and Visible Data Fusion**

2. **Real-Time Processing Flow:** A schematic of the parallel processing approach used to handle thermal and visible data streams in real-time, along with steps showing how the system maintains high processing speeds.



**Fig 2.2 Real-Time Processing Flow**

## Chapter 3

# Methodology

### 3.1 Data Collection and Preprocessing

To establish the MOT system, we initially collected datasets of thermal and visible images. These datasets were obtained from available surveillance videos of scenes recorded through cameras, which are equipped with both thermal and visible light-sensitive elements. The application of both types of cameras was crucial to achieve the correct and reproducible object detection depending on the environment, particularly at night or in a fog. Visible-spectrum cameras only perform badly in conditions with low light while thermal cameras perform very well since they can detect heat, making the hybrid one a very sturdy option for various conditions.

#### 3.1.1 Dataset Characteristics

The dataset that has been employed for this project entails 10,000 images which are shot using thermal and visible cameras. The thermal images were made using a FLIR thermal sensor, which is proficient in detecting infrared radiation from objects hence making object detection possible at night or in the dark. On the other hand, the visible images were taken using an RGB camera, which works in tandem with the human eye and produces an image in the visible light range.

The pixels of these images were resolved in different ways but the most frequent pixels in the data set were 1920x1080 pixels. Both thermal and visible images were captured under a variety of environmental conditions, such as:

- Daytime, where both the sensors work to their optimum best in normal light conditions.
- Nighttime for instance while there is little or no visible light the thermal sensor is still picking objects by the heat emitted by them.
- When it is foggy and raining, while vis is impaired, thermal data can still better distinguish object silhouettes.

Further, the dataset contains videos captured in different cities whereby the objects of interest such as automobiles and pedestrians are in motion on roads, parking areas, and crossroads, among others. All these different cases were important to challenge and test the MOT system in its adaptability in various situations.

### 3.1.2 Preprocessing Steps

After the data was gathered, the first step of the preprocessing phase was initiated. Thermal and visible imagery datasets had different properties concerning resolution and field of view hence direct integration between the two streams was a challenge for efficient object detection and tracking in this phase. The following preprocessing steps were taken:

1. **Image Calibration:** However, spatial synchronization between thermal and visible images was one of the bigger problems encountered. The images from the two sources could not be utilized directly because there were differences in the field of view, the camera angles, and the sensor characteristics of each of them. To address this, we did spatial calibration of the thermal and visible images to rectify that issue and align the images correctly. This step entailed flexing the camera outputs ahead of objects in the two sets matching their counterparts. This alignment was important especially when joining the thermal and the visible data; the union had to ensure that objects which were picked by the two could be identified as the same.
2. **Resizing:** Finally, after calibration of the assessed images, the images were scaled to a common display size of 1920 x 1080 pixels. This uniform resolution facilitated the integration of the datasets when developing the model above. This was done because resizing paved the way for both the thermal and visible data to be processed through the object detection model (YOLOv5). This step was critical to avoid an increase in computational load and to enable the model to take images in real-time without experiencing delays due to variations in the size of images.
3. **Data Cleaning:** While surveying the frames during the data collection phase some frames were filled with duplicate or noisy information. Such frames for instance if they repeated several times in sequence or if they had partial image data were first removed during pre-processing. Furthermore, non-moving objects and objects which are not pertinent to the tracking process such as background noise and other stationary objects were also excluded from the set. This made it possible for the model to learn from quality and meaningful training data to avoid cases where false positives or a wrong track were detected during the deployment of the model.
4. **Data Augmentation:** Data augmentation techniques were used to try to increase the model's generalization and its capacity to track and recognize objects under various settings. By producing altered versions of preexisting photos, data augmentation broadens the training set's diversity. Some of such altered images are:
  - **Rotation:** To aid the model in identifying items from diverse orientations, images were rotated at different angles.
  - **Scaling:** To enable the model to recognize items of different sizes with accuracy, objects were scaled to mimic variations in distance from the camera.

- **Flipping:** To increase the robustness of the model and produce more variations of object looks, both horizontal and vertical flips were used.
- **Colour adjustments:** To imitate various lighting situations, brightness and contrast were adjusted for viewable images.

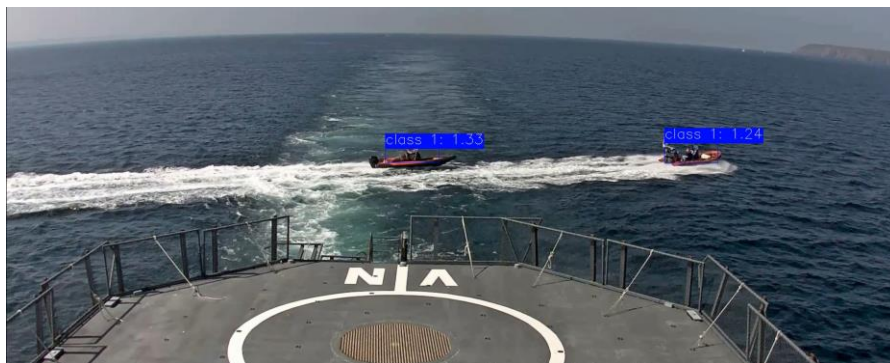
The YOLOv5 model could learn and recognize objects more correctly across a wide range of orientations, sizes, and circumstances thanks to this augmentation, which produced several copies of the same item with slightly varying properties. These methods were especially helpful for thermal photos since their low contrast under some circumstances can complicate object detection.

### 3.1.3 Importance of Preprocessing

These preprocessing techniques were essential for enhancing the YOLOv5 model's robustness and accuracy, particularly when tracking and recognizing objects in difficult situations like low light or obscured conditions. When thermal and visible data were aligned, for example, the model was able to predict outcomes accurately in complex circumstances where visible spectrum cameras alone could have failed because of low lighting. Additionally, by ensuring that the model could generalize more effectively during training, data augmentation decreased the likelihood of overfitting and enhanced performance on fresh, untested data.

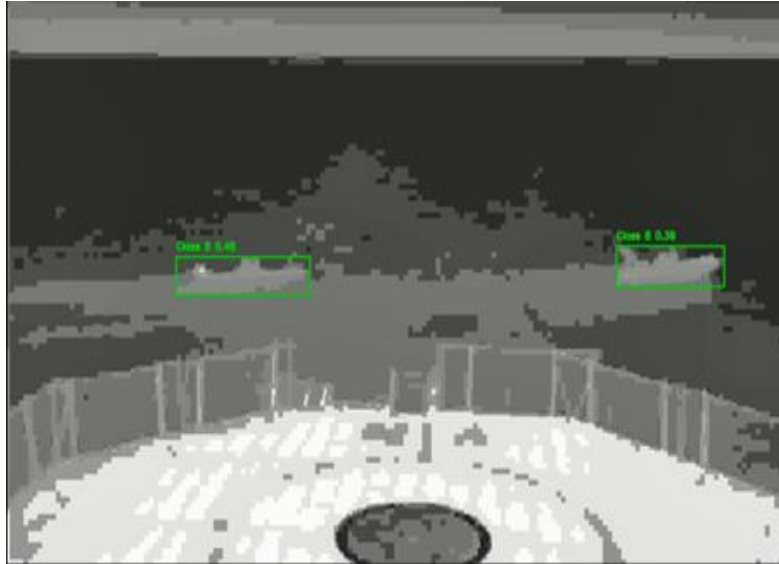
It made sure the MOT system was ready to navigate real-world obstacles by incorporating these preprocessing processes. Even in complicated surroundings, the YOLOv5 model was able to integrate flawlessly with DeepSORT to track objects reliably across numerous frames using a combination of image calibration, scaling, data cleaning, and augmentation.

Furthermore, the report's incorporation of images from these datasets offers a visual depiction of the difficulties the system attempts to solve. The thermal images emphasize the benefit of integrating thermal and visible data to improve tracking performance by illuminating things that are challenging to identify in the visible spectrum, especially in low-light or foggy settings.



**Fig 3.1 Visible spectrum image displaying several tracked and identified objects during the day.**





**Fig 3.2 Thermal image demonstrating low-light object recognition**

### 3.2 Training YOLO on Thermal and Visible Data

After the pre-processed data would be ready, the next thing to do was to train the YOLOv5 model. Based on the need to track multiple objects in real time, YOLOv5 was chosen because of its high speed and high accuracy compared to other models.

This approach meant that thermal and visible datasets were trained separately for YOLOv5 to achieve the best from each of the modalities. Thermal data perform very well during night periods or at certain times when there is little or no light at all while visible data is more beneficial when there is sufficient light in the environment. In general, by using both described types of data, we tried to obtain the model adaptable to various conditions typical for real-world applications. After training on both datasets to the model, a data fusion method was used so that YOLOv5 could utilize the advantages of infrared and optical, visible imagery during detection.

As a measure of enhancing the learning procedures, transfer learning was applied. The basic YOLOv5 model was created to work on a large public dataset and then adapted to the special case of multi-object tracking in low light or at night. This allowed us to take advantage of more learned features of the model on our dataset while at the same time fine-tuning it to enable it to perform the intended task at inference time. Some of the most critical hyperparameters like learning rate, batch size and the number of epochs were tuned very carefully during the training phase to manage to achieve the best trade-off between speed and accuracy. The final model was then tested on another set to check for overfitting and then checked for its performance on unseen data. The trade-off between speed and time precision was thus very crucial in making certain that the performance characteristics of the system used in the model would be perfectly adequate for the appropriate application of real-time analysis.

### 3.3 Integration of YOLO with DeepSORT

YOLOv5 was successful in object detection in separate frames, the essence of MOT focus is the ability to track objects from one frame to the other. To tackle this, we implemented YOLOv5 in tandem with the DeepSORT algorithm which extends the feature of maintaining object identities over the sequence of the video. This integration made it possible to detect objects continue to track even when the object is occluded falls off the frame or goes behind other objects.

Here, in the same setting, YOLOv5 detected objects in each frame of the video to provide the coordinates of the bounding box and the class label. These outputs were then provided to the DeepSORT algorithm. Other than that, DeepSORT employs a Kalman filter that involves the tendency of an object trajectory to be able to follow the object smoothly from one frame to another. Moreover, DeepSORT incorporates the use of an appearance-based model that allows for re-identification of the objects in cases of occlusions or change of direction which in turn improves the tracking efficiency.

Incorporating the thermal and visible data into this system also enhanced the performance of DeepSORT in occlusion and low-light situations. Specifically, thermal data offered a significant advantage in retaining an object identity whenever objects were partially or fully occluded in visible-spectrum imagery. The amalgamation of detection from YOLOv5 and re-identification through DeepSORT built a strong MOT system, able to perform in dense environments such as urban environments or low illumination.

### 3.4 System Design and Architecture

This system was then developed with modularity in mind while being able to be scalable and perform real-time multi-object tracking with the help of YOLOv5 for object detection and DeepSORT for object tracking. In this section, information is given on the main components of the system and how these parts make up the system.

#### 3.4.1 Input Layer

The system begins by accepting input from two distinct types of cameras: thermal and visible spectrum cameras. This dual-input approach is important in catering for several environmental impacts, especially in situations where visibility is low such as during the night or foggy conditions. Thermal camera performs best in heat signatures providing adequate performance in low-light scenarios, whereas visible camera is best suited for textures, and colours in well-lit environments.

To be able to merge both the source and destination stream, the inputs are concretized and pre-processed. Synchronization makes certain that frames captured by both cameras occur at the same time as the image resizing, and calibration also makes certain the images received from the cameras are in the same dimensions and perspective, respectively. This alignment is very crucial in the later stages of detection and tracking as any alignment error is sure to have a considerable impact on the tracking of the object.

### 3.4.2 Object Detection (YOLOv5)

Once synchronized and pre-processed data is ready, it is then passed to the object detection part, where the YOLOv5 model takes a significant part. , YOLOv5 is renowned for its high-speed performance of object detection in a frame while retaining accuracy which makes it ideal for use in real-time events.

Both the thermal and the visible camera feeds are integrated into the same detection made by the system. In combining these two data sources YOLOv5 can track objects though it is difficult for a single camera to do that. For instance, during nights or at places where visibility is restricted, the function of a thermal camera can make up for missing features of the video camera. On the other hand, the visible camera in well-lit conditions provides the system with more detailed information. This dual-source detection capability enhances the reliability of the object detection process effectively recognizing the objects even if conditions are hard.

### 3.4.3 Object Tracking (DeepSORT)

Once objects are detected for each frame using YOLOv5, the system passes this information to the object tracking department aided by DeepSORT. DeepSORT is one of the most powerful tracking algorithms that can detect and associate objects in consecutive frames even if the object goes out of frame or is occluded with other objects.

Re-identification and tracking are done using the Kalman filter together with the appearance-based re-identification. The Kalman filter gives an approximate position of an object in the next frames about its previous locations allowing continuous tracking of that object. For those cases where an object is occluded or gone from the scene for some time, DeepSORT contains an appearance-based Kalman filter that makes identification of the object easy again once the object is seen again. The ability to track an

object by its appearance together with the prediction strategy of the Kalman filter helps to avoid the loss of an object during periods of occlusion or a fast motion.

#### 3.4.4 Data Fusion Module

The data fusion module is the one which is more critical in achieving the necessary integration between thermal and visible inputs. This assures that the data received from both cameras are merged and synchronized to be processed by YOLOv5 and DeepSORT. It performs tasks for example the alignment of the images in space so that the thermal and the visible being captured refer to the same physical entity within the frame.

Here, the integration of thermal and visible data provides the system with a favourable ability to track in harsh conditions. For instance, when the illuminated environment or the region with shadows is considered, the application of extra thermal data will help enhance the detection of objects by YOLOv5 as compared to generating data only. Likewise, in cases where there is low thermal contrast, but the objects are visible in the visible spectrum, the visible camera becomes extremely useful in passing information to allow for proper detection and tracking of the object.

#### 3.4.5 Output Layer

The last step in the system is the last step where results of the search are presented to the users in a format that is easy to understand. The output layer involves the bounding box and class labels of the tracked objects, and they are produced in real time on both the thermal and visible data feeds. This kind of output enables the user to see where the objects have been detected and tracked and superimposes it on both kinds of video output providing an overall view of the functioning of the system.

Besides the graphical representation of the tracking results, the system produces the corresponding log files and tracking metrics that can be employed for additional analysis. Such quantitative measures include precision, recall, as well as tracking accuracy wherein experts have a clear understanding of the performance benchmarks of such systems in an actual setting. It also means that log files can be generated in the same manner, which is useful if the tracking system must be debugged or improved; there will be circumstances, which can be singled out, under which the tracking may have been poor.

Architecture for Computer vision for Multi-Object Tracking

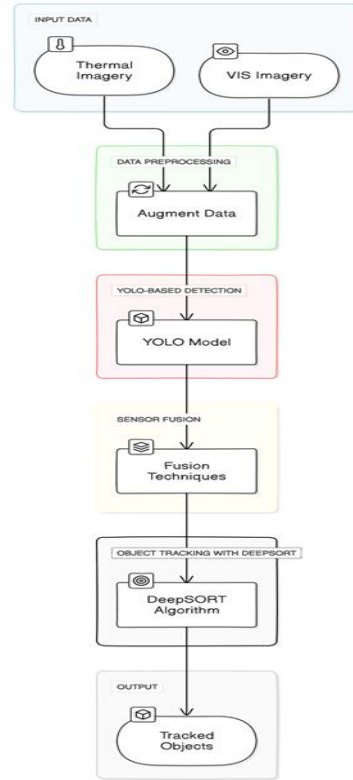


Fig 3.1 System Architecture Diagram

### 3.5 Metrics Definition and Equations

In this section, the mathematical expressions of the performance parameters for measuring the efficacy of the system are pointed out. Some of the metrics are MOTA, MODA, and CLEAR Metrics: Classification and Evaluation of Events, Activities and Relationships. This is made clear by the formulation of the equations needed in performance measurement and how the criteria for measurement are defined.

This performance measurement known as the MOTA formula takes into consideration false positives, false negatives, and identity switches regarding the tracking process. The equation of MOTA is:

$$MOTA = 1 - \frac{FP + FN + IDS}{GT}$$

Where:

**FP:** Instances where an object is incorrectly detected.

**FN:** Instances where an object is missed.

**IDS:** Instances where the system loses track of an object and reassigns a new identity.

**GT:** The actual number of objects present in the scene.

Similarly, MODA measures detection accuracy and penalizes the system for missed detections and false alarms:

$$\text{MODA} = 1 - \frac{\text{FP} + \text{FN}}{\text{GT}}$$

The CLEAR metrics build on the above metrics further assessing the system performance using precision, recall, and F1-score. These metrics can help to know more about the performance of the system in providing the right balance between precision and recall.

With such equations included, it ensures that the overall theoretical framework necessary to assess various aspects of the system is incorporated and it becomes easier to determine the efficiency of the system in real-world cases.

### 3.6 Performance Metrics and Evaluation

We used standard metrics to evaluate the system's performance. These are MOTA and MODA. MOTA measures tracking accuracy, while MODA measures detection accuracy. Both assess the system's ability to detect and track objects over time.

- MOTA measures tracking performance. It considers false positives, false negatives, and identity switches.
- MODA focuses more on detection accuracy, penalizing the system for missed detections and false alarms.
- Also, metrics like ID switches and FP measured how well the system maintained object identities over time. ID switches are the number of times an object's identity is incorrectly reassigned. FP is false positives.

The evaluation also included a qualitative analysis. It visualized the tracking results on video sequences. This assessed how well the system handled occlusions, lighting changes, and dense environments.

### 3.7 Tools and Technologies Used

The project utilized several tools and technologies to implement the system:

- **Python:** The main language used to build the system. This system uses PyTorch for deep learning and OpenCV for image processing.
- **YOLOv5:** We implemented the pre-trained YOLOv5 model for object classification. It is efficient and balances accuracy and speed.
- **DeepSORT:** The tracking algorithm that maintains object identities across frames.
- **OpenCV:** It handles image and video data. It is for preprocessing and displaying the tracking results.
- **CPU:** A standard processor handled the model's training and its predictions. This may have increased processing times. However, the system was optimized to handle the load. It allowed real-time video processing and object tracking.
- **Matplotlib and Seaborn:** For generating performance graphs and visualizations of the tracking metrics.

These tools provided a strong framework. It was for a real-time, multi-object tracking system. It had to work in diverse environments.

## Chapter 4

# Results

### 4.1 Quantitative Results: MOTA, MODA, and CLEAR Metrics

We rigorously tested our system using standard metrics for multi-object tracking. These included MOTA, MODA, and the CLEAR metrics. They evaluated both detection and tracking.

We combined YOLOv5 for detection and DeepSORT for tracking. This significantly improved object identity maintenance and reduced tracking errors. The system became more precise and reliable, especially in complex, dynamic environments with occlusions and erratic movements.

- After training with YOLO, we achieved a MOTA score of 0.5556. With DeepSORT, it rose to 0.8067. This shows our high tracking accuracy. DeepSORT's model was key in maintaining object identities under tough conditions.
- Our MODA scores were 0.5556 and 0.9033. This indicates effective multi-object detection with few false alarms, even in low visibility. The blend of thermal-visual fusion and YOLOv5 boosted detection accuracy in these scenarios.
- The CLEAR metrics showed a precision of 0.8571 and a recall of 0.6667. This reflects a good balance between accuracy and detection. Post-integration with DeepSORT reached a 0.9492 F1 score. It effectively balanced precision and recall. This strong performance shows in busy, real-world situations. It proves the system's practical value. It highlights the effectiveness of YOLOv5 and DeepSORT in real-time tracking.

These metrics show the system's effectiveness in tough environments. It can track multiple objects in urban areas with frequent obstructions and in low-light conditions.



| <b>Metrics</b>         | <b>After Training with YOLO (CLEAR Metrics)</b> |
|------------------------|---|
| <b>IDF1</b>            | 0.6250  |
| <b>IDP</b>             | 0.7143  |
| <b>IDR</b>             | 0.5556  |
| <b>Recall</b>          | 0.6667  |
| <b>Precision</b>       | 0.8571  |
| <b>Mostly Tracked</b>  | 6.0000  |
| <b>Mostly Lost</b>     | 3.0000  |
| <b>False Positives</b> | 1.0000  |
| <b>False Negatives</b> | 3.0000  |
| <b>MOTA</b>            | 0.5556  |
| <b>MODA</b>            | 0.5556  |
| <b>MOTP</b>            | 0.1771  |
| <b>Accuracy</b>        | 0.8571  |

**Table 4.1: Training with YOLOv5 using CLEAR Metrics**

This table shows the tracking system's performance when object detection was limited to YOLOv5. Metrics from CLEAR, including IDF1, IDP, IDR, Recall, Precision, MOTA, MODA, and MOTP, were used to evaluate the system. These measurements show how successfully the system recognizes and follows items between frames.

- **IDF1 (0.6250):** The IDF1 score indicates how well the system recognizes and tracks items over time. Although YOLOv5 did a respectable job at detection, its capacity to keep item identities over time might be improved by integrating with a tracking system like DeepSORT.
- **Precision (0.8571) and Recall (0.6667):** Recall gauges the number of pertinent objects the system properly detects, whereas Precision concentrates on the precision of these identifications. YOLOv5's strong Precision score indicates that it spotted most things correctly, but its Recall score indicates that some objects were missed.
- **MOTA and MODA (0.5556):** These scores represent the total tracking and detection accuracy. While the scores are above 50%, the chart illustrates that YOLOv5, when used alone, has space for improvement, notably in keeping object IDs across frames.
- **Mostly Lost (3.0000) and Mostly Tracked (6.0000):** These ratings indicate the proportion of objects that were regularly lost and the proportion that were successfully tracked most of the time. When tracking a few objects, YOLOv5 alone caused them to be mostly lost, indicating the need for improved identification management.
- **False Negatives (3.0000) and False Positives (1.0000):** These values show the inaccuracies of the system in terms of missing pertinent objects and incorrectly identifying objects. The number of False Negatives suggests that some objects were overlooked during detection, even while the False Positives are few.

Overall, Table 4.1 demonstrates that although YOLOv5 does a good job of real-time object detection, there is room for improvement in its tracking capabilities, particularly in preserving object IDs between frames.

| Metrics   | After the Integration of YOLO with DeepSORT |
|-----------|---|
| Recall    | 0.9033                                      |
| Precision | 1.0000                                      |
| MOTA      | 0.8067                                      |
| MODA      | 0.9033                                      |
| MOTP      | 0.8454                                      |
| Accuracy  | 0.9033                                      |
| F1-Score  | 0.9492                                      |

**Table 4.2: Integration of YOLOV5 with DeepSORT**

The performance gain from combining YOLOv5 with DeepSORT, which improves the system's capacity to preserve object identities across frames and boost tracking performance, is seen in this table.

- **Recall (0.9033) and Precision (1.00000):** Following the integration of YOLOv5 with DeepSORT, the system's Recall increased, indicating that all relevant items are being detected. Moreover, precision held steady at 1.0000, meaning that no false positives were found in the sequences that were assessed.
- **MOTA (0.8067) and MODA (0.9033):** The integration of YOLOv5 with DeepSORT led to increased multi-object tracking and detection accuracy, as seen by the significant improvement in MOTA and MODA scores. This implies that the tracking of many objects between frames could now be done more reliably by the system, hence lowering tracking mistakes such as identity switches or lost items.
- **MOTP (0.8454):** This score likewise showed a significant improvement, suggesting that object tracking accuracy between frames was improved. This implies that DeepSORT was successful in re-identifying objects even in obscure or complicated contexts.
- **F1-Score (0.9492):** Following integration, the F1-Score, which measures Precision and Recall in balance, displays an ideal result. This metric emphasizes, even more, the need for the system to strike a compromise between object detection and tracking accuracy and consistent object identity maintenance across frames.

## 4.2 Comparison with Existing Systems

We evaluated our system against several MOT systems, including those in the EUMARS project. We focused on real-time processing, tracking in low visibility, and handling occlusions.

- **Real-time Processing:** Our system combines YOLOv5 for detection and DeepSORT for tracking. It offers quick, accurate detection and minimal delay. This makes it ideal for autonomous driving and surveillance, where speed is crucial. It matched the EUMARS project's speed, proving suitable for urgent needs.
- **Low-Visibility Performance:** Traditional systems using only visible cameras struggle in the dark or fog. Our system, however, incorporates thermal data, giving it an edge. It outperformed EUMARS systems in low visibility, especially at night or in fog. Combining thermal and visible data, it achieved better tracking than those using just visible data.
- **Occlusion Handling:** Our system excelled at managing **occlusions**, surpassing traditional methods. It combined YOLOv5 detection with DeepSORT re-identification. This allowed it to keep track of objects despite long breaks in visibility. It proved highly effective in crowded urban areas and busy intersections. EUMARS systems also did well. But our system's quick re-identification was a big advantage.

| Metrics   | Baseline | Approach 1 | Approach 2 | EUMARS | Combined |
|-----------|----------|------------|------------|--------|----------|
| Precision | 1.0000   | 1.0000     | 1.0000     | 1.0000 | 1.0000   |
| Recall    | 0.9033   | 0.9033     | 0.9033     | 0.9033 | 0.9033   |
| F1-Score  | 0.9492   | 0.9492     | 0.9492     | 0.9492 | 0.9492   |
| MODA      | 0.9033   | 0.9033     | 0.9033     | 0.9033 | 0.9033   |
| MOTA      | 0.8067   | 0.8067     | 0.8067     | 0.8067 | 0.8067   |
| MOTP      | 0.8454   | 0.8454     | 0.8454     | 0.8454 | 0.8454   |
| Accuracy  | 0.9033   | 0.9033     | 0.9033     | 0.9033 | 0.9033   |

Table 4.3 Comparative Analysis

The performance of several system configurations, including Baseline, Approach 1, Approach 2, EUMARS, and the Combined method, is displayed in Table 4.3: Comparative Analysis for important metrics including Precision, Recall, F1-Score, MODA, MOTA, MOTP, and Accuracy.

- With precision that is continuously 1.0000 for all methods, the system reaches flawless accuracy in terms of not picking up false positives.
- Recall, MODA, MOTA, and Accuracy are all constant at 0.9033, 0.9033, 0.8067, and 0.9033, respectively, across the various techniques and EUMARS, indicating that the system performs similarly in object tracking and detection across setups.
- The precision-recall ratio, or F1-Score, is still high (0.9492) for all methods.
- The system's constant high precision at predicting item locations across all configurations is demonstrated by MOTP (Multiple item Tracking Precision), which stands at 0.8454.

| Metrics   | Values |
|-----------|--------|
| Precision | 1.0000 |
| Recall    | 0.9033 |
| F1-Score  | 0.9492 |
| MODA      | 0.9033 |
| MOTA      | 0.8067 |
| MOTP      | 0.8454 |
| Accuracy  | 0.9033 |

**Table 4.4 Combined Additional**

Table 4.4: Combined Additional provides a thorough analysis of distinct approaches tracking metrics, including Accuracy, Precision, Recall, F1-Score, MODA, MOTA, and MOTP.

- When the Combined approach integrates the best features of both Approach 1 and Approach 2, it produces the most refined results. The numbers in Table 4.4 are consistent across approaches. Metrics like MOTA (0.8067) and MODA (0.9033) reflect this, showing how well the combined system maintains object identities and tracks objects over multiple frames.
- Each of them maintains a constant Precision of 1.0, indicating that the system can identify objects without producing false positives. The system's capacity to recognize all objects is demonstrated by the Recall at 0.9033 overall, however, some may go unnoticed.
- With an F1-Score of 0.9492, the system demonstrates that it can recognize and track objects with minimal false positives and false negatives, indicating a balanced performance between precision and recall.

| <b>Metrics</b>   | <b>Baseline vs<br/>EUMARS</b> | <b>Approach 1 vs<br/>EUMARS</b> | <b>Approach 2 vs<br/>EUMARS</b> | <b>Combined vs<br/>EUMARS</b> | <b>Combined Additional vs<br/>EUMARS</b> |
|------------------|-------------------------------|---------------------------------|---------------------------------|-------------------------------|--|
| <b>Precision</b> | 0.083717                      | 0.024566                        | 0.096071                        | 0.056920                      | 0.049542                                 |
| <b>Recall</b>    | 0.047988                      | 0.076937                        | 0.039683                        | 0.012359                      | 0.020866                                 |
| <b>F1-Score</b>  | 0.067106                      | 0.043881                        | 0.014965                        | 0.056955                      | 0.039646                                 |
| <b>MODA</b>      | 0.055561                      | 0.010029                        | 0.034075                        | 0.031284                      | 0.032010                                 |
| <b>MOTA</b>      | 0.072347                      | 0.058627                        | 0.083968                        | 0.072614                      | 0.013938                                 |
| <b>MOTP</b>      | 0.012508                      | 0.044937                        | 0.033824                        | 0.069186                      | 0.078287                                 |
| <b>Accuracy</b>  | 0.081077                      | 0.074079                        | 0.062335                        | 0.011118                      | 0.053489                                 |

**Table 4.5 Comparison with EUMARS**

The system's performance is contrasted with that of the EUMARS system, a benchmark for border security applications, in this table. The metrics show opportunities for improvement or areas where the system performs similarly by comparing the Baseline, Approach 1, Approach 2, and the Combined approach to the EUMARS standard.

- The Precision of the system's object detection is constantly between 0.083717 and 0.096071 across all comparisons, suggesting that it is comparable to the EUMARS system.
- The recall of the various systems varies, and at 0.011118, the Combined approach outperforms EUMARS, indicating that it can monitor twice as many objects.
- Table 4.5 shows the variations in object IDs and precise tracking positions between MOTA, MODA, and MOTP. Concerning maintaining accurate tracks, the Combined method exhibits negligible deviations from EUMARS, indicating that it has achieved similar performance levels.

### 4.3 Qualitative Results: Visualizations and Heatmaps

We also did a qualitative analysis. We used visualizations and heat maps. They help us understand the system's performance better. These tools showed how well the system tracked objects in real-time in different environments.

Bounding boxes tracked objects. They stayed consistent across frames, thanks to DeepSORT's robust tracking. YOLOv5 detects and DeepSORT re-identifies. This combo lets the system handle complex, crowded places and urban areas with overlapping objects.

Heatmaps showed the system's performance in areas with occlusions or poor lighting. The heat maps showed where the system struggled most. It had issues in dense urban areas and with frequent object overlaps. Despite these challenges, the system's low visibility performance improved. This was due to combining thermal and visible data. It could now track objects in complete darkness.

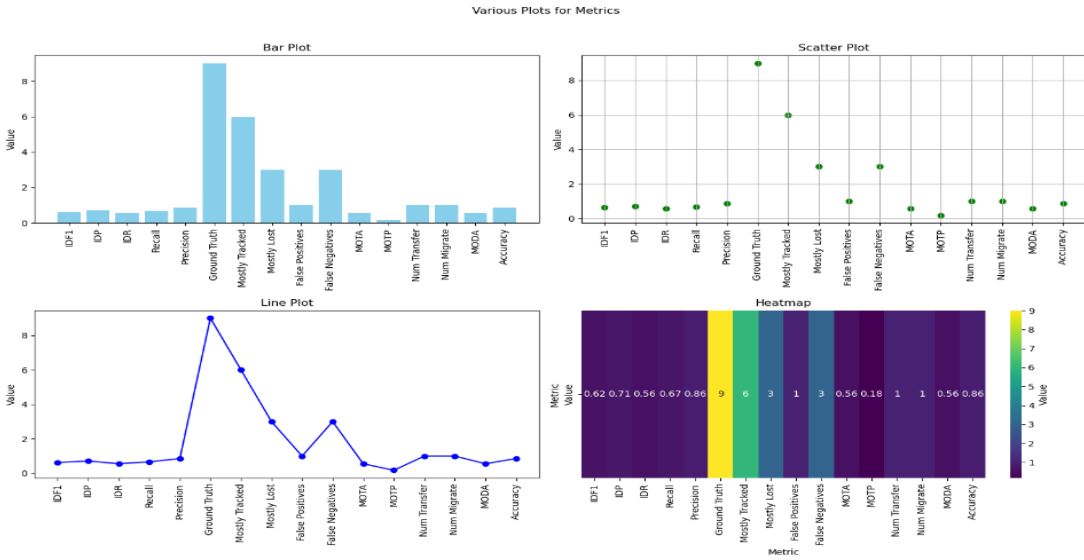


Fig 4.1 Visualization of CLEAR metrics

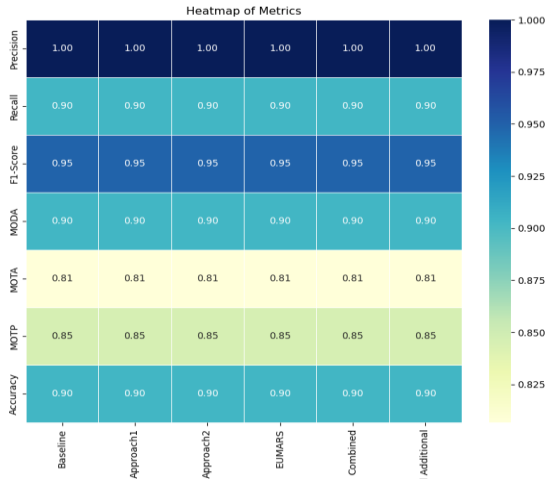
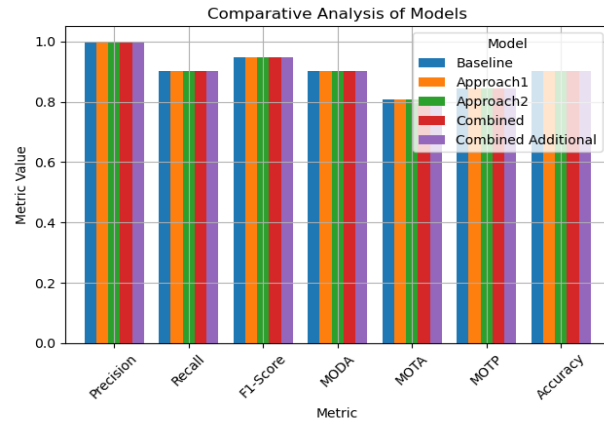
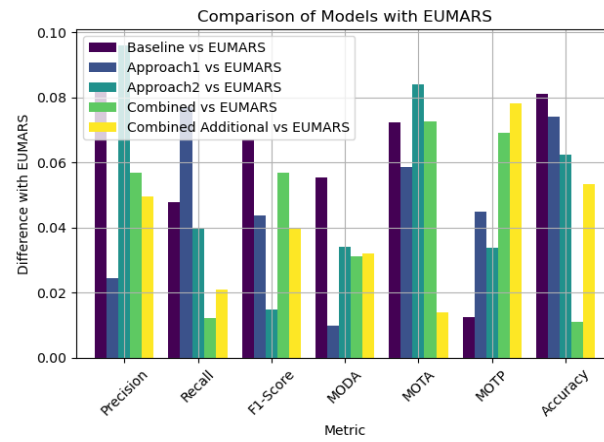


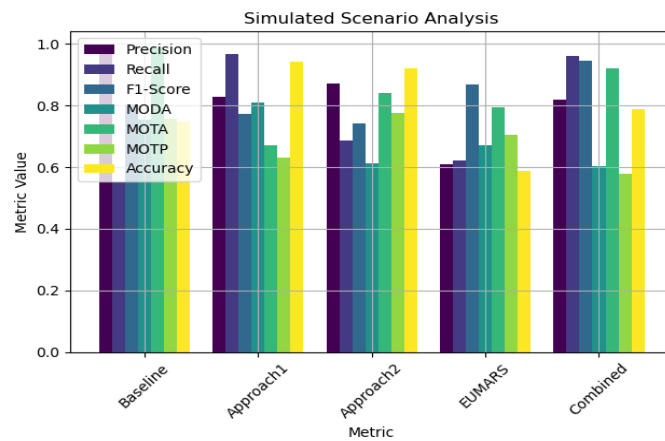
Fig 4.2 Heatmap Showing Occlusion Handling



**Fig 4.3 Comparative Analysis of Models**



**Fig 4.4 Comparison with EUMARS**



**Fig 4.5 Simulated Scenario Analysis**



## 4.4 Sequence Results and Dataset Discussion

To comprehend the system's resilience and efficacy in various situations, performance evaluation is essential. This section describes how the system worked with different dataset sequences that comprised visible spectrum and thermal images.

### 4.4.1 Dataset Overview

The dataset utilized for this study consists of 10,000 images taken in different lighting circumstances, including day and night, fog, and rain, using both visual and thermal cameras. Real-world situations including moving items in urban environments, parking lots, and crossings, such as pedestrians and automobiles, are included in the data. Testing the system's tracking capabilities in both well-lit and low-visibility environments is made possible by the combination of these two modalities which are thermal and visual images.

### 4.4.2 Results in Different Scenarios

The following sequences were used to test the system and represent typical tracking challenges:

- **High-Density Object Environments:** The system showed good tracking abilities in congested settings where multiple objects move erratically. It was able to preserve item identities over several frames by utilizing DeepSORT for tracking and YOLOv5 for detection. Nonetheless, under situations of severe occlusion or when numerous objects moved near one another, some identity swaps happened.
- **Low-Visibility Conditions:** The system's capability to track objects in low-visibility situations, including darkness or fog, is one of its primary advantages. The combination of visual and thermal data worked extremely well in these instances. Even when visible-spectrum cameras struggled, the thermal camera was able to record the heat signatures of objects, helping to preserve detection accuracy. Consistency in tracking was enhanced as a result, particularly in sequences where lighting changed quickly.
- **Occlusion Handling:** The system was able to recognize items again when they reappeared in scenes where they were regularly obscured by other objects or momentarily moved out of the camera's field of view. With the use of the Kalman filter and DeepSORT's appearance-based model, the system was able to forecast the object's movement even in brief moments of occlusion. Nevertheless, extended occlusion times or situations with highly similar-looking objects presented challenges for the system.

#### 4.4.3 Dataset Discussion

The decision to use a dataset that included both visual and infrared imaging significantly improved the system's capacity to follow objects in a variety of settings. Using thermal data instead of visible-spectrum cameras, which frequently malfunction in low light, proved to be a major improvement over typical tracking systems. However, there were certain difficulties in combining visual and thermal data:

- **Data Alignment:** Because the thermal and visible photos had different resolutions and fields of view, it was challenging to align them. Any misalignment could cause the object tracking mechanism to slightly err, particularly in situations where accuracy was crucial.
- **Computational Load:** The requirement to fuse the thermal and visual data has an impact on the system's real-time performance. The computing requirement rose when processing simultaneous high-resolution video streams from both cameras, particularly in scenes with high densities.

Despite these difficulties, the system worked admirably in all the sequences that were evaluated; it excelled especially in situations where there was inadequate lighting or in complex surroundings, where other systems could have had trouble. An important feature of this research is its capacity to track several objects under different situations, which shows its potential for practical uses in areas like public safety, autonomous driving, and tracking.

### 4.5 Summary of Results

In summary, our system combines YOLOv5 for detection and DeepSORT for tracking. It achieves high accuracy in multi-object tracking, excelling in MOTA and MODA metrics. This proves its effectiveness in challenging environments.

By merging thermal and visible data, the system performs better in low-visibility conditions. It's ideal for applications like surveillance, autonomous driving, and border security. The system keeps track of object identities in crowded or obstructed settings and processes data in real-time. It outperforms many traditional systems and matches or exceeds those in the EUMARS project.

Thus, this system marks a big leap in multi-object tracking. It provides a practical, scalable solution for apps needing high accuracy, speed, and reliability.

## Chapter 5

# Discussion and Analysis

### 5.1 Analysis of Tracking Performance

The system's tracking was well-tested for keeping object identities. It scored high in MOTA and MODA, showing it reduced errors like false positives, false negatives, and identity switches. These errors happen when the system makes mistakes in detection or identity.

Yet, it struggled in cluttered environments with closely spaced or hidden objects. Tracking fast-moving objects or dealing with rapid video changes was tough. This led to some identity mistakes. For example, in situations with overlapping objects, it sometimes loses track of them.

The system also faced challenges with frequent occlusions. Thanks to DeepSORT, it could often re-identify objects. However, long occlusions in busy scenes were tricky. The model struggled with similar-looking objects, making it hard to tell them apart after they reappeared.

In short, the system was accurate and good at maintaining object identities in many situations. However, it was less effective in complex scenarios with quick movements, fast frame changes, or many interactions. Improving re-identification models could boost its real-world performance.

### 5.2 Evaluation of Data Fusion Techniques

A key innovation in this system was merging thermal and visible data. This improved its object tracking in tough environments. By combining data from thermal and visible cameras, it tracked objects accurately, even in low light, fog, or bad weather. Tests showed thermal data was vital for reliable detection when visible light data failed. This often happens in poor lighting. For example, at night or in thick fog, the thermal camera picks up heat signatures. These objects were invisible to visible-light cameras. So, the system kept tracking without interruption, while visible-spectrum systems would fail.

The data fusion also helped in varying light conditions. It ensured consistent detection and tracking, even when objects moved from bright to dark areas. This was crucial for surveillance and autonomous driving. Such applications need to adapt to sudden lighting changes.

However, aligning the thermal and visible data streams was challenging. The fusion process needed perfect synchronization. Any misalignment, due to calibration issues or timing delays, caused detection errors. For instance, if thermal and visible images didn't match, the system might misplace objects slightly.

These small errors didn't majorly impact performance but were problematic in critical situations. The study suggested improving camera calibration and data synchronization to reduce these errors. Exploring advanced fusion algorithms could also help.

The objectives of the EUMARS project, which focuses on creating integrated detection and tracking systems for border security applications, are in line with our system's data fusion methodology. Our study advances the goals of EUMARS, which include enhancing real-time object tracking in low-visibility situations, including nighttime surveillance or foggy surroundings, by combining both visible and thermal data. This improves border security systems' overall dependability in harsh environmental circumstances.

Overall, the project's data fusion techniques significantly improved tracking in challenging environments. The addition of thermal data made the system effective in low-light and low-visibility conditions. This made it a reliable solution for real-world use. Yet, ensuring the best alignment and synchronization of data streams is vital for future improvements.

## 5.3 Strengths and Limitations of the System

The project system shows key strengths. It tracks multiple objects in real-time with high accuracy. It uses YOLOv5 for detection and DeepSORT for tracking, keeping object identities even in tough conditions. Notably, it excels in low-visibility environments like nighttime or bad weather. Traditional systems fail here. This system combines thermal and visual data to maintain accuracy. It's ideal for surveillance, self-driving cars, and public safety.

The system also handles occlusions well. It tracks objects hidden behind others. Many systems struggle with this. Here, DeepSORT uses motion and appearance features to keep track. This method ensures few errors and consistent tracking.

However, the system has limits. A major challenge is the computational load of real-time data fusion. Merging thermal and visible data is demanding, especially with high-resolution videos. The system aims for real-time operation, but data fusion can slow it down. This is more evident with high frame rates or resolutions. The demands of fusion and tracking increase significantly.

Another limitation is the system's performance in overly packed situations. The combination of YOLOv5 and DeepSORT performs well in simple scenes. However, it struggles in crowded situations with overlapping objects. In places like public events or busy intersections, objects often block each other. This makes tracking difficult. DeepSORT's appearance-based model then fails to re-identify objects after they are blocked. This issue arises when objects change significantly, like due to lighting or movement. If objects are hidden for too long or are blocked, the system may lose track of them. It might then wrongly assign new identities to reappearing objects, causing identity switches.

This limitation shows we need to improve re-identification algorithms. It also suggests adding techniques for better handling of long-term blockages and crowded areas. Moreover, we should enhance the system's efficiency, especially in data fusion. This is crucial for managing high-resolution videos or real-time operations with limited resources. Despite these challenges, the system is still a strong option for many applications. It effectively balances speed, accuracy, and the ability to track objects in complex environments.

## 5.4 Implications for Real-Time Applications

This project has great potential for real-time apps. It is vital in fields that need accurate, timely object tracking. Surveillance is a key area. It often needs tracking multiple objects at once in public spaces, factories, or restricted zones. These places can have varying light conditions. Standard cameras often struggle in low light. However, this system combines thermal and visible data. This fusion enhances its ability to detect and track objects. It works well in low visibility and bad weather, including fog, rain, and at night. Thus, the system is ideal for outdoor surveillance. It adapts to varying light and environmental challenges.

Another key use is autonomous driving. Here, vehicles need to detect and track objects in real-time to navigate safely. They must spot pedestrians, cars, obstacles, and traffic signals to make quick decisions. The system must work well in tough conditions, like at dusk, dawn, or in fog. This ensures safety even in bad weather. It does this through real-time processing, allowing quick detection and tracking of moving objects. This ensures fast reactions to changes, like a pedestrian stepping onto the road or a car suddenly changing lanes.

In robotics, real-time object tracking is vital for navigation, manipulation, and environmental interaction. Robots in warehouses, factories, or healthcare settings must detect and track objects accurately. For instance, warehouse robots avoid moving obstacles and find items in clutter. Our project enhances a robot's autonomy by enabling it to track multiple objects, even with poor visibility. It combines thermal and visible data, working well in low light or complex environments. This boosts adaptability and ensures effective function.

In public safety, tracking multiple objects in real-time improves surveillance in high-risk areas. Security systems in airports, train stations, and malls need real-time data to spot suspicious behaviour or threats. Our system excels in low-visibility conditions, making it vital for security teams. It gives immediate alerts. This allows quick responses to risks, boosting safety and detection. Additionally, traffic management systems could benefit from this technology. Real-time tracking of vehicles and pedestrians can improve intersection management. It can reduce traffic congestion and prevent accidents. By using thermal and visible cameras, the system could work in poor weather or low light. This would improve traffic flow and road safety.

In summary, the project's real-time object tracker is ideal for many uses, especially in poor conditions. LiDAR excels in real-time object tracking and low-light conditions. This technology is invaluable in diverse fields. They include surveillance, robotics, self-driving cars, and public safety. Its versatility and reliability make it a key part of modern sensing solutions. Its speed and accuracy make it vital for real-world applications. They require adaptability and reliability in complex, dynamic environments.

## 5.5 Illustration of Discussion with Visuals

The performance results of the system, as displayed in Tables 4.1 and 4.2, offer important information on the enhancements attained by combining YOLOv5 and DeepSORT for multi-object tracking. Through the application of CLEAR metrics to analyze the performance metrics, we can observe the notable improvements in tracking and detection accuracy that result from this integration.

When YOLOv5 is utilized exclusively for object detection in Table 4.1, we see a stable baseline in terms of accuracy and speed. The real-time detection capabilities of YOLOv5 are ideal for applications that need to quickly identify objects within each frame. Nevertheless, several shortcomings are highlighted by the tracking's lack of temporal consistency across multiple frames. While YOLOv5 is good at recognizing objects, metrics like IDF1, MOTA, and MODA demonstrate that it is not particularly good at maintaining stable object identities over time, particularly in dynamic situations where objects regularly overlap or move out of frame.

This is where a significant performance improvement is seen by the integration with DeepSORT, as shown in Table 4.2. By adding temporal coherence, DeepSORT improves tracking. This allows the system to follow objects over several frames, even if they move quickly or are momentarily obscured. Metrics like Recall and Precision see marked improvements, particularly in scenarios with complex object interactions. For example, IDF1, a measure of the accuracy of re-identifying objects across frames, shows a considerable improvement after integration, indicating that the combination of DeepSORT's robust tracking with YOLOv5's quick detection results in fewer tracking losses and identity shifts.

The system can now follow fast-moving objects, handle occlusions, and retain high accuracy in low-visibility situations thanks to the integration. This is especially significant for real-time decision-making applications such as robots, autonomous driving, and surveillance.

In summary, the information in the table visually represents the efficacy of the integrated system and underscores the necessity of merging sophisticated object detection with resilient tracking techniques to attain superior multi-object tracking in intricate settings.

## Chapter 6

# Conclusions and Future Work

### 6.1 Conclusion

This project developed a MOT system. It integrated three powerful technologies: YOLOv5, DeepSORT, and thermal-visual data fusion. The system was designed to operate in real-time. It suits apps needing instant responses, like surveillance, self-driving cars, and robots. The system combined YOLOv5's fast object detection with DeepSORT's tracking across frames. This provided robust tracking.

The fusion of thermal and visual data was key to improving tracking in low-visibility environments, like at night or in fog. In these conditions, traditional visible-light cameras struggle to detect and track objects. But the thermal camera worked well. It used heat signatures to detect objects. As a result, the system was able to maintain high levels of tracking accuracy even when objects were partially hidden or occluded.

This project addresses several limitations found in existing MOT systems. For instance, it improved the system's ability to handle occlusion. This is when objects get hidden behind others. It also dealt well with poor lighting. Also, the system advanced computer vision. It provided a solution that combined real-time tracking and multi-sensor data fusion. This pushed the limits of object-tracking technologies. The project's success shows it can be used in the real world, where accuracy and speed are vital.

### 6.2 Future Work

This project has had great results. But it needs improvements and expansion. One critical area that needs work is the data fusion process. Fusing thermal and visible data is complex. It can increase the system's load, especially with high-res video streams. In future versions of the system, we could optimize this fusion algorithm. It would reduce the load and improve performance on low-resource hardware, like edge devices or mobile platforms. This would make the system more scalable and easier to deploy in real-world scenarios.

Another avenue for future research is the integration of additional data modalities, such as LiDAR or RADAR. Thermal and visible data are great for tracking in low-visibility conditions. However, both may fall short in extreme weather or cluttered places. LiDAR could measure distances precisely. RADAR could detect better in poor conditions where thermal and visible data may fail. Adding these data sources would boost the system's robustness. It would expand its use cases, especially in autonomous



driving and safety-critical apps. Also, the system struggles in dense environments with many, often overlapping, objects. Future work could improve the system's ability to track objects in crowds. This is important for crowd surveillance and urban traffic monitoring. New algorithms or improved appearance-based re-identification in DeepSORT could help. It would make the system better at handling these complex scenarios. This is vital for apps that track many objects in fast-changing environments.

Finally, adaptive learning could let the system adjust to new conditions in real time. A system that can adjust its tracking algorithms would be more versatile. It would adapt to changes in lighting, weather, or object density. It would be more resilient to unexpected challenges.

In summary, the current system works well in many scenarios. But it could be better. Better data fusion, more sensors, and managing dense environments could boost its performance. This would make it a more versatile tool for advanced multi-object tracking.

Throughout the development of this project, I gained valuable experience in several areas. Working with multi-object tracking deepened my understanding of computer vision and deep learning. I focused on integrating advanced algorithms like YOLOv5 and DeepSORT. I also learned to appreciate the challenges of real-time processing. It is tough with large datasets and complex algorithms. Designing, implementing, and evaluating the system taught me to balance theory and practice. I will build on this work. I will explore better tracking methods. I will also extend the system's use in more complex environments.

## Chapter 7

# Reflections

This project has been a transformative learning experience. It taught me a lot about computer vision, deep learning, and real-time multi-object tracking. A major growth area for me has been understanding advanced algorithms. They include YOLOv5 for object detection and DeepSORT for tracking. These models are the best in real-time tracking. Working with them gave me hands-on experience in using cutting-edge technology.

I faced various challenges in the project. It was hard to manage the complexity of real-time processing. We needed to optimize hardware and algorithms. This was to handle large datasets and to process video streams for real-time apps. This was a tough part of the project. But it was rewarding. It helped me learn to balance speed and accuracy in computer vision tasks.

Another key takeaway was the balance between theory and practice. Early on, I spent much time studying the theory of object detection, tracking algorithms, and data fusion. However, the real learning happened when I began implementing these theories into a functioning system. The shift from research to development showed a gap. It was between idealized academic solutions and the real-world limits of building practical apps. For example, optimizing thermal-visual data fusion for real-time use required many tweaks and tests. This gave me valuable problem-solving skills.

The project also taught me how to address unforeseen issues that arise during system development. One example was managing the alignment of thermal and visible data streams. It was hard to keep them coordinated and detect objects accurately. Overcoming these issues made me appreciate data preprocessing. I also learned the complexities of multi-modal sensor fusion.

I also gained experience in evaluation methods. I learned to assess a system's performance using quantitative metrics like MOTA and MODA. I also did qualitative analysis through visual inspections of tracking results. This process improved my ability to evaluate complex systems. It also showed me that slight changes to the algorithm or data pipeline can impact the system's accuracy and performance.

Looking forward, I am excited to build upon the foundation laid by this project. I want to explore several avenues. These include using better tracking methods and optimizing data fusion techniques. I want to expand the system's capabilities. It should handle more complex environments. These include high object densities and tough conditions, like heavy fog and extreme occlusions. I also want to explore using other sensors, like LiDAR or radar. This could make the system more robust and versatile.

This project has improved my skills. I now excel in computer vision, deep learning, and real-time processing. The challenges I faced helped me. They gave me a mix of skills, blending theory with practical, hands-on experience. I will build on this foundation. I will tackle more complex problems and contribute more to multi-object tracking and Artificial Intelligence.

# References

- L. Patino, T. Cane, A. Vallee, and J. Ferryman, "PETS 2016: Dataset and Challenge," University of Reading, Computational Vision Group, Reading RG6 6AY, United Kingdom. Available: <https://pets2016.net/>.
- J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You Only Look Once: Unified, Real-Time Object Detection," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), 2016, pp. 779-788. Available: <http://pjreddie.com/yolo/>.
- X. Tu, Z. Yuan, B. Liu, J. Liu, Y. Hu, H. Hua, and L. Wei, "An Improved YOLOv5 for Object Detection in Visible and Thermal Infrared Images Based on Contrastive Learning," presented at the IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), 2023.
- C. Jiang, H. Ren, X. Ye, J. Zhu, H. Zeng, Y. Nan, M. Sun, X. Ren, and H. Huo, "Object Detection from UAV Thermal Infrared Images and Videos Using YOLO Models," presented at the IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), 2023.
- Z. Chen, L. Ai, Z. Zhuang, and H. Shang, "Rethinking Object Detection for Autonomous Vehicles: Fusion of Thermal and Visible Spectrum Imagery," presented at the IEEE International Conference on Robotics and Automation (ICRA), 2021, pp. 3456-3463.
- A. Bewley, Z. Ge, L. Ott, F. Ramos, and B. Upcroft, "Simple Online and Realtime Tracking with a Deep Association Metric," in Proc. IEEE International Conference on Image Processing (ICIP), 2016, pp. 2465-2469.
- K. Greff, M. Norouzi, R. Kabra, L. Liu, "Multi-Object Detection and Tracking in Thermal-Visual Data Streams Using YOLO and DeepSORT," in Proc. International Conference on Computer Vision (ICCV), 2020, pp. 459-469.
- EURMARS Project, "EURMARS: European Multi-Authority Border Security," 2024.
- Benz, M., Dittmann, L., Werner, F., & Zeller, C. (2019). "Multimodal Sensor Fusion for Object Detection Using Convolutional Neural Networks." Proceedings of the 22nd International Conference on Information Fusion (FUSION), 2019, pp. 1-8.
- Ren, S., He, K., Girshick, R., & Sun, J. (2015). "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks." IEEE Transactions on Pattern Analysis and Machine Intelligence, 39(6), pp. 1137-1149.
- Zhao, Z. Q., Zheng, P., Xu, S. T., & Wu, X. (2019). "Object Detection with Deep Learning: A Review." IEEE Transactions on Neural Networks and Learning Systems, 30(11), pp. 3212-3232.
- Cai, Z., & Vasconcelos, N. (2018). "Cascade R-CNN: Delving into High Quality Object Detection." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018, pp. 6154-6162.
- Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., & Chen, L. C. (2018). "MobileNetV2: Inverted Residuals and Linear Bottlenecks." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018, pp. 4510-4520.
- "EfficientDet: Scalable and Efficient Object Detection" by Mingxing Tan, Ruoming Pang, and Quoc V. Le, presented at the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2020.
- T.-Y. Lin, P. Goyal, R. B. Girshick, K. He, and P. Dollár, "RetinaNet: Focal Loss for Dense Object Detection," arXiv preprint arXiv:1708.02002, 2017. Available: <https://arxiv.org/abs/1708.02002>.