

D4.2.1 Person Detection and Localization

Due date: **31/12/2023**
Submission Date: **10/01/2024**
Revision Date: **21/02/2025**

Start date of project: **01/07/2023**

Duration: **36 months**

Lead organisation for this deliverable: **Carnegie Mellon University Africa**

Responsible Person: **Yohannes Haile**

Revision: **1.1**

Project funded by the African Engineering and Technology Network (Afretec) Inclusive Digital Transformation Research Grant Programme	
Dissemination Level	
PU	Public
PP	Restricted to other programme participants (including Afretec Administration)
RE	Restricted to a group specified by the consortium (including Afretec Administration)
CO	Confidential, only for members of the consortium (including Afretec Administration)

Executive Summary

Deliverable D4.2.1 focuses on the development of a ROS node that detects and localizes people under various conditions. This deliverable includes the implementation of a ROS node called `personDetection`, accompanied by a comprehensive report documenting the development process, refinement of requirements, and a detailed specification of the node's functional characteristics. Additionally, it provides a user manual with clear instructions on building and launching the ROS node. The design of the interface covers input, output, and control data, with suitable data structures and code that adhere to the software engineering standards established by the project. The functionality of the `personDetection` node is thoroughly tested and validated using various test cases, including scenarios with different lighting conditions, occlusions, and varying distances between the robot and the user. The node is also tested on the Pepper robot to confirm its reliability and real-time performance, ensuring it meets the intended objectives effectively.

Contents

1	Introduction	4
2	Requirements Definition	5
3	Module Specifications	6
4	Module Design	7
5	Implementation	10
6	Running the Person Detection Node	14
7	Unit Test	15
	References	17
	Principal Contributors	18
	Document History	19

1 Introduction

This document outlines the development and implementation of a ROS node for person detection, localization. The primary goal of this node is to enhance the interaction capabilities of the Pepper robot, allowing it to identify and track people within its environment, which serves as the behavior controller to initiate interaction as specified in behavior specification.

The deliverable includes a detailed report documenting the complete software development life cycle for the person detection and localization module. The requirements definition process is thoroughly covered, ensuring that all functional necessities are carefully aligned with the project's objectives. This section also highlights any identified misalignment or challenges that may arise during the development process and how they are addressed. The module design section provides an in-depth description of the person detection localization functionality, covering critical aspects such as input, output, and control data.

The operation of the module is guided by parameters defined in a configuration file, which is structured as a list of key-value pairs in the `personDetectionConfiguration.json` file. This configuration allows for flexible and scalable customization of the module's behavior to suit various operating conditions and requirements. Furthermore, the document emphasizes the importance of robust design principles to ensure the module's reliability and performance in real-time applications.

2 Requirements Definition

The person detection and localization module is designed to meet the following requirements, ensuring seamless integration with the Pepper robot and its ROS-based ecosystem. The key requirements for the module are as follows:

Person Detection

- Detect people in the robot's field of view using an RGB image as input.
- Detect and localize all people in the field of view when multiple people are present.
- Localize person by determining their position in the image and drawing bounding boxes.
- Identify the centroid coordinates of each bounding box.
- Determine the depth (distance from Pepper robot) of all the people detected.

Person Labeling and Consistency

- Assign unique labels to detected person (e.g., "Person 1").
- Maintain consistent labeling of the same person across consecutive image frames, provided the spatial displacement is within a defined tolerance.
- Reassign new labels to reappearing person if not detected for a configurable number of images.

Configurable Parameters

- Allow customization through a configuration file (`personDetectionConfiguration.json`)

Input/Output Specifications

- **Input:** RGB-D image from a robot camera or external camera.
- **Output:** Annotated images with bounding boxes and labeled person records published to the `/personDetection/data` topic.

Verbose Mode

- Provide optional diagnostic output and visual debugging through an OpenCV window.

Misalignment of the module

Due to the poor camera present on the robot camera it necessitated the use of external camera. Hence, the depth information provided by the Pepper's camera is low quality hence the depth information (distance of the people from the camera) isn't accurate. In addition, this module doesn't support the simulator.

3 Module Specifications

The person detection module, implemented as a ROS node named `personDetection`, is designed to detect people within Pepper robot's field of view and determine their location in the image frame of reference. The module provides labeled, color-coded bounding boxes around each detected person and tracks these label across successive images for coherent detection.

The inputs for this module is an RGB image from the robot's camera or external camera (Intel RealSense D435i), the depth iamge from the robot's depth sensors or an external RGB-D camera (Intel Realsense D435i).

The outputs for this module is annotated RGB image with bounding boxes around detected people and an array of record is published with the following message `personDetection/data` topic:

- Person label representing as number
- 3D image coordinates of the bounding box centroid

YOLO is a deep learning framework optimized for real-time object detection, and when tailored for person detection for accurately identifying human figures in diverse environments. Once persons are detected, the system integrates SORT (Simple Online and Realtime Tracking) for tracking across frames. SORT uses techniques like Kalman filtering and the Hungarian algorithm to maintain consistent identities for each detected individual, even as they move or occlude one another. This combination between YOLO's rapid detection and SORT's efficient tracking provides a robust solution for real-time monitoring for Pepper robot.[1]

If `verbosemode` is set to True in the configuration file, an OpenCV window will display the detected person's bounding box. Each detected person will also be assigned a unique label. This provides real-time visualization and tracking for person detection.

A unit test is developed to cover various scenarios, including multiple people, partial occlusion, variable lighting conditions, and label reassignment when the person disappear. The tests will be conducted using a driver-stub test platform, which utilizes recorded color and depth images stored in the data folder. Additionally, the unit tests can be executed directly on the physical robot to validate real-world performance.

4 Module Design

Image Input

The primary input for the ROS node will be the Intel RealSense camera mounted on top of Pepper's head. As an alternative, the Pepper camera can also be used by configuring the camera parameter in the configuration file. However, as noted in section 2 the depth camera has very low quality. The Intel RealSense camera provides both RGB and depth images at various resolutions and frame rates, which can be customized through the launch file parameters. The table below outlines the available resolution and frame rate configurations for the Intel RealSense camera.

Format	Resolution	Frame Rate (FPS)	Comment
Z [16 bits]	1280x720	6, 15, 30	Depth
	848x480	6, 15, 30, 60, 90	
	640x480	6, 15, 30, 60, 90	
	640x360	6, 15, 30, 60, 90	
	480x270	6, 15, 30, 60, 90	
	424x240	6, 15, 30, 60, 90	
YUY2 [16 bits]	1920x1080	6, 15, 30	Color Stream from RGB camera (Camera D415 & D435/D435i)
	1280x720	6, 15, 30	
	960x540	6, 15, 30, 60	
	848x480	6, 15, 30, 60	
	640x480	6, 15, 30, 60	
	640x360	6, 15, 30, 60	
	424x240	6, 15, 30, 60	
	320x240	6, 30, 60	
	320x180	6, 30, 60	

Table 1: Stream Configurations for Depth and Color for Intel RealSense D435i. See the official datasheet: [Intel RealSense D400 Series Datasheet](#).

Algorithms

YOLO (You look only Once)

The YOLO algorithm is an deep learning method for real-time object detection. It uses a single convolutional neural network to predict multiple bounding boxes and class probabilities for each object in the frame. The model divides the image into a grid and makes predictions for each cell simultaneously, allowing it to detect multiple objects quickly and accurately. YOLO's architecture significantly reduces computation time, making it suitable for real-time applications. During post-processing, overlapping boxes are filtered using non-max suppression to ensure precision and reduce false positives. When employed for person detection specifically, YOLO efficiently isolates human figures in complex environments, enabling fast and reliable identification even in dynamic scenes. In our implementation, we are using the YOLOV8s model, a state-of-the-art variant optimized for enhanced speed and accuracy. YOLOV8s retains the core strengths of the original YOLO framework while offering improvements in efficiency and detection performance. YOLOV8 comes with five variants based on the number of parameter nano(n), small(s), medium(m), large(l), and extralarge(x). We are employ-

ing the small(s) variant which is designed for resource-constrained environments and delivers faster inference.

SORT (Simple Online and Realtime Tracker)

The SORT algorithm is a lightweight multi-object tracking method that combines Kalman filtering for motion prediction and the Hungarian algorithm for data association. The process begins with detecting a person using an Intel RealSense camera mounted on Pepper's head. YOLO is employed for head detection, after which the Kalman filter predicts the motion of detected objects in subsequent frames. To associate new detections with existing tracks, SORT utilizes the Hungarian algorithm with Intersection over Union (IoU) as the matching criterion. Once matches are found, the Kalman filter updates its state with the latest information. Tracks that do not find a match are marked as lost and eventually deleted, while new detections initiate new tracks.[2]

Algorithm 1 SORT Algorithm

Require: Detected bounding boxes B_t at time t , tracked objects O_{t-1} from time $t - 1$

Ensure: Updated object IDs and bounding boxes O_t

```
1: if  $O_{t-1}$  is empty then
2:   for all bounding box  $b \in B_t$  do
3:     Register  $b$  as a new object with a unique ID
4:   end for
5: else
6:   Predict new positions of tracked objects using the Kalman filter
7:   Compute the cost matrix  $D$  using Intersection over Union (IoU) between  $B_t$  and predicted
   objects
8:   Solve the assignment problem using the Hungarian algorithm
9:   for all matched pairs  $(o, b)$  do
10:    Update object  $o$  with new bounding box  $b$ 
11:    Reset disappearance counter for  $o$ 
12:  end for
13:  for all unmatched objects in  $O_{t-1}$  do
14:    Increment disappearance counter
15:    if counter exceeds threshold then
16:      Deregress the object
17:    end if
18:  end for
19:  for all unmatched bounding boxes in  $B_t$  do
20:    Register bounding box  $b$  as a new object with a unique ID
21:  end for
22: end if
23: return updated objects  $O_t$ 
```

Below is a diagram illustrating the complete process of SORT tracking.

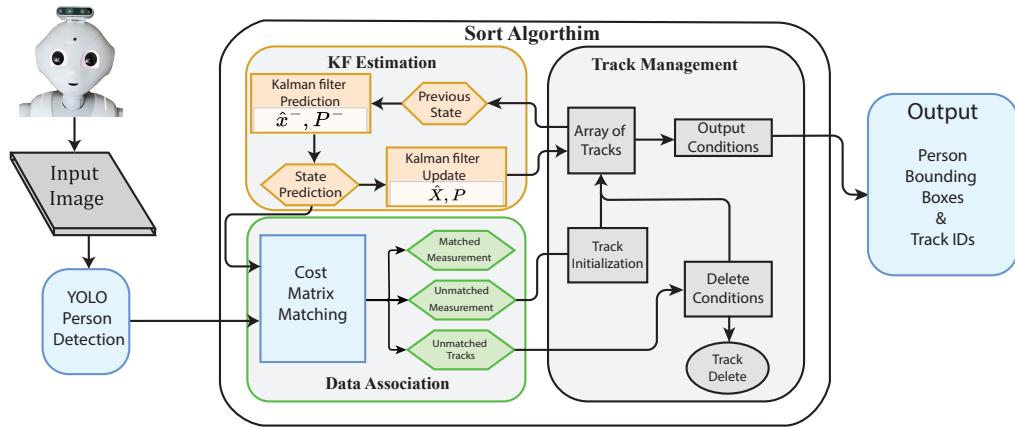


Figure 1: SORT Diagram

5 Implementation

File Organization and Purpose

The source code for conducting person detection and localization is structured into three primary components: `person_detection_application`, `person_detection_implementation`, and `person_detection_tracking`. The `person_detection_implementation` component encapsulates all the essential functionality required for executing YOLO based person detection. The `person_detection_tracking` component, on the other hand, manages tracking functionality by employing SORT (Simple Online and Realtime Tracking). Additionally, the person detection system is equipped with the capability to process various files critical for testing, such as configuration files, input files, and topic files. Meanwhile, the `person_detection_application` component serves as the entry point, invoking the main functions to run the person detection node and executing the functions defined within `person_detection_implementation`.

Here is the file structure of the person detection package:

```
cssr_system
└── person_detection
    ├── config
    │   └── person_detection_configuration.json
    ├── data
    │   └── pepper_topics.dat
    ├── models
    │   └── person_detection_yolov8s.onnx
    ├── msg
    │   └── msg_file.msg
    ├── src
    │   ├── person_detection_application.py
    │   ├── person_detection_implementation.py
    │   └── person_detection_tracking.py
    ├── CSSR4AfricaLogo.svg
    ├── person_detection_requirements.txt
    ├── README.md
    └── CMakeLists.txt
    └── Package.xml
```

Figure 2: File structure of the person detection system

UML Diagram for the Person Detection and Localization Module

The UML diagram provides a clear structural representation of the Persons Detection and Localization Module, illustrating the relationships between its core components. It highlights inheritance, where `PersonDetectionNode` serves as the base class, extended by `YOLOv8` for specialized person detection. Associations between tracking components such as `Sort`, and `TrackerUtils` emphasize how detected people are tracked using Kalman filtering.

Below is the UML diagram of person_detection_implementation.py

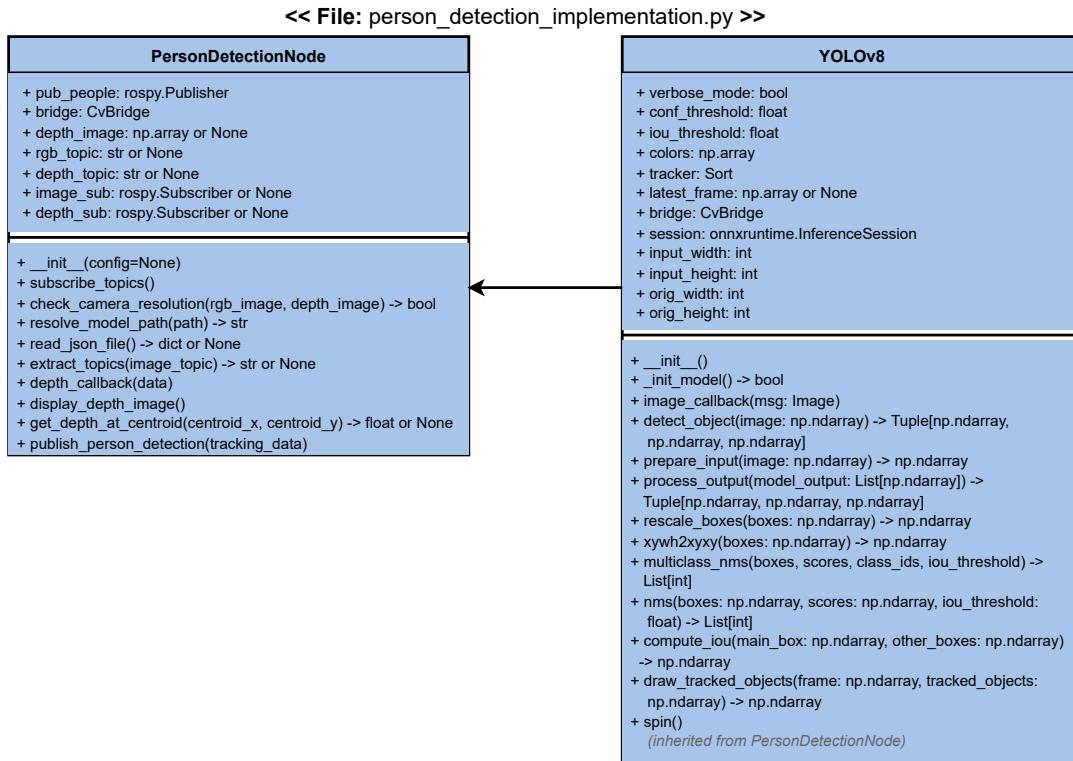


Figure 3: Person detection implementation UML

Below is the UML diagram of person_detection_tracking.py

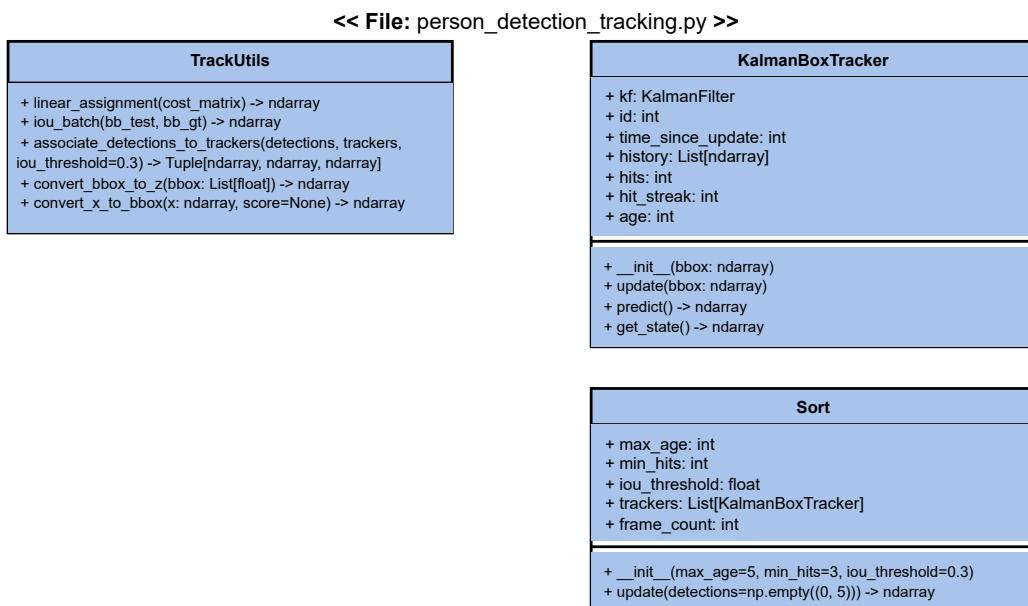


Figure 4: Person detection tracking UML

Configuration File

The operation of the person detection node is determined by the contents of the configuration file that contains a list of key-value pairs as shown on the table below. The configuration file is named `person-detection-configuration.json`.

Key	Value	Description
YOLO_confidence	<number>	Specifies the confidence threshold for the YOLO person detection algorithm.
sort_max_disappeared	<number>	Specifies the maximum number of frames an object can disappear for Sort tracker before being removed.
sort_min_hits	<number>	Specifies the minimum number of consecutive hits required for Sort tracker initialization.
sort_iou_threshold	<number>	Specifies the Intersection over Union (IoU) threshold for Sort tracker associations.
verbose_mode	true, false	Specifies whether diagnostic data is to be printed to the terminal and diagnostic images are to be displayed in OpenCV windows.

Table 2: Configuration file key-value pairs for the persons detection node.

Input File

There is no input file the person detection node.

Output File

There is no output file the person detection node. The node using OpenCV to display the detected people with bounding boxes and labels.

Models

The person detection node uses one model for detecting people. The model is stored in the `models` directory. The model is:

Model	Description
<code>person-detection-yolov8s.onnx</code>	YOLO-based person detection model.

Table 3: Models used by the person detection node.

Topics File

For the test, a selected list of the topics for the robot is stored in the topics file. The topic files are written in the `.dat` file format. The data file is written in key-value pairs where the key is the camera and the value is the topic. The topics file for the robot is named `pepper-topics.dat`.

Launch File

The launch file `person_detection_launch_robot.launch` is designed to initialize either Pepper's front camera or the Intel RealSense camera based on the specified configuration. It declares several parameters that can be customized to match your network settings and camera choice:

- `pepper_robot_ip`: specifies the IP address of the Pepper robot (default: `172.29.111.230`).
- `pepper_robot_port`: specifies the communication port for Pepper (default: `9559`).
- `network_interface`: specifies the network interface name (default: `wlp0s20f3`).
- `namespace`: sets the ROS namespace for the `naoqi` driver (default: `naoqi_driver`).
- `camera`: selects the camera source; set to `pepper` for Pepper's front camera or `realsense` for the Intel RealSense camera (default: `realsense`).
- `unit_test`: determines whether to run unit tests (default: `False`).

The file sets the parameter `/personDetection/camera` to the chosen camera and conditionally launches the corresponding nodes. If the `camera` parameter is set to `pepper`, the launch file starts the `naoqi_driver` node using the provided IP, port, network interface, and namespace. Conversely, if `camera` is set to `realsense`, it includes the RealSense camera launch file with specified parameters for image resolution, frame rate, and depth alignment. Users can adjust these default values to suit their specific hardware configurations.

Topics Subscribed

The person detection node subscribes to the following topics:

Camera	Topic Name	Message Type
RealSenseCameraRGB	<code>/camera/color/image_raw</code>	<code>sensor_msgs/Image</code>
RealSenseCameraDepth	<code>/camera/aligned_depth_to_color/image_raw</code>	<code>sensor_msgs/Image</code>
PepperFrontCamera	<code>/naoqi_driver/camera/front/image_raw</code>	<code>sensor_msgs/Image</code>
PepperDepthCamera	<code>/naoqi_driver/camera/depth/image_raw</code>	<code>sensor_msgs/Image</code>

Table 4: Topics subscribed by the person detection node.

Topics Published

The person detection node publishes the following topics:

Topic Name	Message Type	Description
<code>/personDetection/data</code>	<code>personDetection/msg_file</code>	An array of records containing person labels and 3D image coordinates of the bounding box.

Table 5: Topics published by the person detection node.

6 Running the Person Detection Node

To run the person detection node, the user must first install the necessary software packages as outlined in [Deliverable 3.3](#). The required packages are listed in the `person_detection_requirements.txt` file. The user can follow the README file in the person detection package to install the required packages. Referring to the implementation section of this deliverable report, the user must set the configuration file to the desired parameters. Using the key-value pair, the user can set the camera, algorithm, confidence threshold, and other parameters. The user can then run the person detection node by executing the following command in the terminal:

```
# Launch either Pepper Camera or RealSense Camera from the launch file
$ roslaunch person_detection person_detection.launch camera:=pepper
# or
$ roslaunch person_detection person_detection.launch camera:=realsense

# Run the person detection node
$ rosrun person_detection person_detection_application.py
```

If the user has set the verbose mode to True in the configuration file, the person detection node will display the detected people with bounding boxes and labels in an OpenCV window.

7 Unit Test

The unit test is designed to validate the person detection node's functionality under various scenarios, including multiple people, occlusions, and varying lighting conditions. The test can be performed using a driver-stub test platform, which utilizes recorded color and depth images stored in the data folder as a rosbag file. The unit test can also be executed directly on the physical robot to validate real-world performance.

The person detection unit test file structure is as follows:

```

unit_test
└── person_detection_test
    ├── config
    │   └── person_detection_test_configuration.json
    ├── data
    │   ├── person_detection_input_singlePerson.bag
    │   ├── person_detection_input_multiplePerson.bag
    │   ├── person_detection_input_personTracking.bag
    │   ├── person_detection_input_occlusion.bag
    │   └── person_detection_input_lighting.bag
    ├── launch
    │   ├── person_detection_launch_robot.launch
    │   └── person_detection_launch_testHarness.launch
    ├── src
    │   ├── person_detection_test_application.py
    │   └── person_detection_test_implementation.py
    ├── CMakeLists.txt
    ├── Package.xml
    └── README.md

```

Figure 5: File structure of the person detection unit test.

The test cases for the person detection node that are going to be evaluated are as follows:

Test Case	Description
Single Person Detection	Verify the person detection node's ability to detect and localize a single person in the image frame, as well as evaluate the distance at which the person is detected.
Multiple Person Detection	Validate the person detection node's capability to detect and localize multiple people in the image frame.
Person Tracking	Test the person detection node's tracking functionality by tracking a person across multiple frames.
Occlusion Handling	Evaluate the person detection node's performance in handling partial occlusions of person.
Lighting Conditions	Test the person detection node's robustness under varying lighting conditions.

Table 6: Test cases for person detection node evaluation.

Configuration File

The configuration file for the person detection unit test is named `person_detection_test_configuration.json` and contains the following key-value pairs:

Key	Value	Description
<code>bag_file</code>	<code><test name></code>	Specifies the ROS bag file used as input for testing.
<code>save_video</code>	true or false	Specifies whether to save the output video of the test.
<code>save_image</code>	true or false	Specifies whether to save individual image frames from the test.
<code>video_duration</code>	<code><number></code>	Specifies the duration (in seconds) for which the video is saved.
<code>image_interval</code>	<code><number></code>	Specifies the time interval (in seconds) at which images are captured and saved.
<code>speaker</code>	true or false	Enables the speaker to announce which test is currently running.
<code>verbose_mode</code>	true or false	Specifies whether detailed logs and diagnostic images are displayed during execution.

Table 7: Configuration file key-value pairs for the person detection test.

Note: Valid values for `bag_file` include: `singlePerson`, `multiplePerson`, `personTracking`, `occlusion`, `lighting`.

References

- [1] Joseph Redmon and Ali Farhadi. Yolo9000: Better, faster, stronger. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7263–7271, 2017.
- [2] Nicolai Wojke, Alex Bewley, and Dietrich Paulus. Simple online and realtime tracking with a deep association metric. *arXiv preprint arXiv:1703.07402*, 2017.

Principal Contributors

The main authors of this deliverable are as follows (in alphabetical order).

Yohannes Haile, Carnegie Mellon University Africa.

Natasha Mutangana, Carnegie Mellon University Africa.

David Vernon, Carnegie Mellon University Africa.

Document History

Version 1.0

First draft.

Natasha Mutangana.

10 January 2024.

Version 1.1

Complete rewrite of the deliverable.

Yohannes Haile.

21 February 2025.